

FFCI: A Framework for Interpretable Automatic Evaluation of Summarization

Fajri Koto
Timothy Baldwin
Jey Han Lau

School of Computing and Information Systems
The University of Melbourne
Victoria 3010, Australia

FFAJRI@STUDENT.UNIMELB.EDU.AU
TB@LDWIN.NET
JEYHAN.LAU@GMAIL.COM

Abstract

In this paper, we propose FFCI, a framework for fine-grained summarization evaluation that comprises four elements: faithfulness (degree of factual consistency with the source), focus (precision of summary content relative to the reference), coverage (recall of summary content relative to the reference), and inter-sentential coherence (document fluency between adjacent sentences). We construct a novel dataset for focus, coverage, and inter-sentential coherence, and develop automatic methods for evaluating each of the four dimensions of FFCI based on cross-comparison of evaluation metrics and model-based evaluation methods, including question answering (QA) approaches, semantic textual similarity (STS), next-sentence prediction (NSP), and scores derived from 19 pre-trained language models. We then apply the developed metrics in evaluating a broad range of summarization models across two datasets, with some surprising findings.

1. Introduction

Remarkable advances in abstractive summarization in recent years have unfortunately not been accompanied by commensurate improvements in automatic evaluation metrics. Most recent studies (Nallapati et al., 2016; See et al., 2017; Gehrmann et al., 2018; Liu & Lapata, 2019; Zhang et al., 2020a; Lewis et al., 2020) continue to rely on ROUGE (Lin, 2004), a lexical-overlap metric that is not capable of detecting paraphrases in abstractive summaries relative to reference summaries, with the only real mainstream alternative being manual evaluation (Hsu et al., 2018; Chen & Bansal, 2018; Hardy & Vlachos, 2018; Celikyilmaz et al., 2018; Krishna & Srinivasan, 2018).

We identify four key dimensions across which to evaluate summaries: (1) *f* faithfulness (Maynez et al., 2020), (2) *f* ocus, (3) *c* overage, and (4) *i* nter-sentential coherence; we label the combined approach “FFCI”. Faithfulness measures the degree of factual consistency (and lack of hallucination) relative to the source document, and is especially important for abstractive methods. The other three aspects are inspired by manual evaluation in previous work (Peyrard & Gurevych, 2018; Hsu et al., 2018; Celikyilmaz et al., 2018; Narayan et al., 2018b; Chen & Bansal, 2018), as summarized by Hardy et al. (2019).

We first revisit recent work on *faithfulness*, and propose a simpler automatic evaluation scheme. Recent work has used question generation (QG) and question answering (QA) to evaluate faithfulness (Wang et al., 2020; Durmus et al., 2020). However, we argue that this

Gold summary : Info-A; Info-B; Info-C
System summary:
Good <i>focus</i> , and Good <i>coverage</i> : Info-A; Info-B; Info-C
Good <i>focus</i> , and Bad <i>coverage</i> : Info-A; Info-A
Bad <i>focus</i> , and Good <i>coverage</i> : Info-A; Info-B; Info-C; Info-D; Info-E
Bad <i>focus</i> , and Bad <i>coverage</i> : Info-D; Info-E; Info-F

Figure 1: Illustration of focus and coverage.

approach is computationally expensive¹ and critically depends on resources that are often unavailable in languages other than English. As an alternative, we extend the experiments of Zhang et al. (2020b) and Durmus et al. (2020) in investigating scores from a broad range of pre-trained language models (computed between summary and article) and find them to be more reliable than QA-based methods.

Secondly, we propose *focus* and *coverage* relative to the reference summary. Both assess semantic equivalence, with *focus* evaluating the proportion of important information in the generated summary (= precision), and *coverage* evaluating the degree of salient information in the reference summary that the generated summary contains (= recall). In Figure 1, we illustrate different scenarios for focus and coverage of system summaries.

Lastly, we address the automatic evaluation of linguistic quality in multi-sentence summaries. Previous work has looked at aspects including readability, fluency, and clarity (Hardy et al., 2019), but we argue that *inter-sentential coherence* is more important for evaluating abstractive summaries for two reasons. First, modern pre-trained language models are highly adept at generating fluent sentences, but global coherence beyond the sentence is not a given. Second, inter-sentential coherence subsumes sub-sentence coherence, as disfluent sentences will break the global discourse coherence.

To summarize, our contributions are: (1) we release an annotated dataset for evaluating focus, coverage, and inter-sentential coherence; (2) for faithfulness, focus and coverage, we benchmark traditional metrics such as ROUGE, METEOR, and BLEU with model-based metrics, including question-answering (QA) methods, semantic textual similarity (STS), FactCC (Kryscinski et al., 2020), and scores from 19 different pre-trained language models; (3) we adapt the next sentence prediction (NSP) for evaluating inter-sentential coherence; and (4) we re-evaluate a broad range of contemporary summarization models over CNN/DailyMail and XSUM based on FFCI, with a number of surprising findings not captured by ROUGE. Data and code used in this paper can be accessed at <https://github.com/fajri91/ffci>.

1. For instance, to evaluate the 11,490 CNN/DailyMail test set (Hermann et al., 2015) requires the generation of roughly 229,800 questions and answers.

2. Related Work

2.1 Aspects on Summarization Evaluation

Automatic evaluations of language generation systems have been based on the comparison of reference and system-generated text. BLEU (Papineni et al., 2002) is a precision-based metric in machine translation task, while ROUGE (Lin, 2004) is the de facto metric for summarization systems (See et al., 2017; Liu & Lapata, 2019; Zhang et al., 2020a). In the other text generation tasks such as caption generation (Xu et al., 2015) and question generation (Du et al., 2017), CIDEr (Vedantam et al., 2015) and SPICE (Anderson et al., 2016) are used to complement BLEU and ROUGE. Recently, pre-trained embedding based evaluation metrics such as BERTSCORE (Zhang et al., 2020b) and MOVERSCORE (Zhao et al., 2019) have also been proposed.

Paper	Automatic					No manual eval	Manual											
	ROUGE	METEOR	BLEU	BERTScore	MoverScore		Faithfulness	Recall	Precision	Relevance	Coherence	Fluency	Relative	Absolute	SCU	reference article	ref+article	Quality control
See et al. (2017)	✓	✓				✓												
Yang et al. (2017)	✓							✓				✓			✓			
Lin et al. (2018)	✓					✓												
Cohan et al. (2018)	✓					✓												
Liao et al. (2018)	✓					✓												
Kedzie et al. (2018)	✓					✓												
Amplayo et al. (2018)	✓					✓			✓	✓	✓					✓		
Jadhav and Rajan (2018)	✓					✓												
Li et al. (2018a)	✓					✓												
Pasunuru and Bansal (2018)	✓	✓				✓												
Cao et al. (2018)	✓					✓												
Sakaue et al. (2018)	✓					✓												
Celikyilmaz et al. (2018)	✓					✓		✓	✓	✓	✓	✓	✓			✓		
Chen and Bansal (2018)	✓	✓				✓			✓	✓	✓	✓	✓			✓		
Guo et al. (2018)	✓	✓				✓			✓	✓	✓	✓	✓			✓		
Hardy and Vlachos (2018)	✓		✓			✓				✓		✓						
Hsu et al. (2018)	✓					✓		✓	✓		✓		✓			✓		✓
Krishna and Srinivasan (2018)	✓					✓		✓			✓		✓			✓		✓
Kryściński et al. (2018)	✓					✓		✓			✓		✓			✓		✓
Li et al. (2018b)	✓					✓	✓					✓						
Narayan et al. (2018a)	✓					✓		✓			✓		✓			✓		
Narayan et al. (2018b)	✓					✓		✓			✓		✓			✓		
Narayan et al. (2018c)	✓					✓		✓			✓		✓			✓		
Peyrard and Gurevych (2018)	✓					✓		✓	✓		✓		✓			✓		
ShafieiBavani et al. (2018)	✓					✓		✓		✓	✓		✓	✓				
Song et al. (2018)	✓					✓	✓	✓			✓		✓			✓		
Hardy et al. (2019)	✓					✓	✓	✓			✓		✓			✓		
Makino et al. (2019)	✓					✓												

Table 1: Evaluation methods used in previous work (split over 3 pages)

Paper	Automatic				No manual eval	Manual										
	ROUGE	METEOR	BLEU	BERTScore MoverScore		Faithfulness	Recall	Precision	Relevance	Coherence	Fluency	Relative	Absolute	SCU	reference article	ref+article
Zhong et al. (2019)	✓				✓											
Peyrard (2019a)							✓	✓				✓	✓			
Fabbri et al. (2019)	✓						✓	✓		✓	✓				✓	
Lev et al. (2019)	✓															
You et al. (2019)	✓				✓											
Isonuma et al. (2019)	✓				✓											
Wang et al. (2019b)	✓						✓	✓		✓		✓			✓	✓
Lebanoff et al. (2019)	✓						✓						✓			
Li et al. (2019a)	✓		✓		✓											
Sharma et al. (2019)	✓				✓											
Wang et al. (2019a)	✓				✓											
Abacha and Demner-Fushman (2019)	✓				✓											
Duan et al. (2019a)	✓				✓						✓					
Zhang et al. (2019)	✓														✓	
Kouris et al. (2019)	✓				✓											
Zhou and Rush (2019)	✓				✓											
Zheng and Lapata (2019)	✓					✓										
Frermann and Klementiev (2019)	✓						✓			✓	✓	✓			✓	
Palaskar et al. (2019)	✓		✓				✓		✓	✓	✓	✓		✓		
Sun and Nenkova (2019)	✓						✓	✓	✓	✓			✓			
Gui et al. (2019)	✓	✓							✓	✓	✓				✓	
Xiao and Carenini (2019)	✓	✓			✓											
Luo et al. (2019)	✓						✓	✓			✓				✓	
Duan et al. (2019b)	✓				✓											
Zhu et al. (2019)	✓						✓	✓		✓	✓	✓				
Liu and Lapata (2019)	✓						✓	✓		✓					✓	
Gao et al. (2019)	✓				✓											
West et al. (2019)	✓							✓		✓	✓					
Shen et al. (2019)	✓					✓				✓					✓	
Parida and Motlicek (2019)	✓		✓		✓											
Grenander et al. (2019)	✓				✓											
Li et al. (2019b)	✓			✓					✓	✓	✓				✓	
Shapira et al. (2019)							✓			✓	✓	✓	✓	✓		
Falke and Gurevych (2019)	✓	✓					✓	✓		✓	✓				✓	
Liu et al. (2019a)	✓						✓	✓		✓	✓				✓	
Ouyang et al. (2019)	✓						✓			✓	✓				✓	
Kim et al. (2019)	✓								✓		✓				✓	
Mendes et al. (2019)	✓						✓		✓						✓	
Koto et al. (2020)	✓			✓			✓	✓								✓
Huang et al. (2020b)	✓				✓											
Xiao and Carenini (2020)	✓				✓											
Lebanoff et al. (2020)	✓				✓											
Xu et al. (2020a)	✓				✓											
Kano et al. (2020)	✓				✓											

Table 1: Evaluation methods used in previous work (split over 3 pages)

Paper	Automatic					No manual eval	Manual												
	ROUGE	METEOR	BLEU	BERTScore	MoverScore		Faithfulness	Recall	Precision	Relevance	Coherence	Fluency	Relative	Absolute	SCU	reference article	ref+article	Quality control	
Gholipour Ghalandari et al. (2020)	✓					✓													
Zhu et al. (2020)	✓				✓			✓	✓				✓						
Gholipour Ghalandari and Ifrim (2020)	✓					✓													
Xu et al. (2020b)	✓								✓				✓				✓		
Gao et al. (2020)									✓				✓	✓					
Sotudeh Gharebagh et al. (2020)	✓						✓	✓				✓			✓				
Maynez et al. (2020)	✓			✓			✓	✓								✓			
Amplayo and Lapata (2020)	✓	✓					✓	✓		✓	✓						✓		
Mao et al. (2020a)	✓	✓					✓	✓				✓							
Xu et al. (2020)	✓	✓					✓			✓	✓		✓						
Schumann et al. (2020)	✓	✓					✓				✓	✓	✓						
Ladhak et al. (2020)	✓	✓	✓				✓	✓				✓	✓	✓					
Durmus et al. (2020)	✓	✓		✓			✓	✓				✓	✓	✓					
Huang et al. (2020c)	✓	✓					✓	✓				✓	✓	✓		✓			
Bražinskis et al. (2020)	✓	✓					✓	✓	✓		✓	✓	✓	✓		✓			
Suhara et al. (2020)	✓	✓					✓	✓	✓		✓	✓	✓	✓		✓			
Li et al. (2020)	✓	✓					✓	✓	✓		✓	✓	✓	✓		✓			
Zhong et al. (2020)	✓	✓				✓													
Huang et al. (2020a)	✓	✓					✓		✓				✓			✓		✓	
Wang et al. (2020a)	✓	✓						✓	✓				✓						
Wang et al. (2020b)	✓	✓				✓					✓								
Mao et al. (2020b)	✓	✓				✓													
Xiao et al. (2020)	✓	✓				✓													
Wu et al. (2020)	✓	✓	✓		✓			✓	✓	✓		✓	✓			✓		✓	
Jia et al. (2020)	✓	✓						✓				✓					✓		
Xu and Lapata (2020)	✓	✓						✓	✓		✓	✓	✓						
Zou et al. (2020)	✓	✓						✓	✓		✓	✓	✓						
Desai et al. (2020)	✓	✓					✓	✓			✓	✓	✓			✓		✓	
Cao et al. (2020)							✓	✓				✓	✓			✓			
Tan et al. (2020)	✓	✓					✓	✓		✓		✓	✓						
Deng et al. (2020)	✓	✓					✓	✓	✓		✓	✓	✓						
Scialom et al. (2020)	✓	✓	✓			✓													
Lu et al. (2020)	✓	✓										✓				✓			
Pilault et al. (2020)	✓	✓					✓	✓		✓	✓	✓	✓					✓	
Bhandari et al. (2020)	✓	✓						✓	✓			✓	✓	✓				✓	
Zhao et al. (2020)	✓	✓						✓			✓	✓	✓					✓	
Cui et al. (2020)	✓	✓				✓						✓	✓				✓		
Lee et al. (2020)	✓	✓					✓		✓			✓	✓			✓			
He et al. (2020)	✓	✓				✓						✓	✓			✓			
Total	106	11	4	5	2	40	18	46	25	12	13	45	26	37	6	7	34	9	7

Table 1: Evaluation methods used in previous work (split over 3 pages)

Paper	Automatic					No manual eval	Manual												
	ROUGE	METEOR	BLEU	BERTScore	MoverScore		Faithfulness	Recall	Precision	Relevance	Coherence	Fluency	Relative	Absolute	SCU	reference	article	ref+article	Quality control
Percentage	95	2	2	2	2	64	10	15	10	5	5	15	10	10	2	2	10	2	2

Table 1: Evaluation methods used in previous work (split over 3 pages)

To comprehensively understand how recent summarization work has performed evaluation, we followed the lead of Hardy et al. (2019) in conducting a survey of 111 summarization papers from major NLP conferences over the period 2017–2020, and group them into automatic and manual evaluation methods in Table 1.² Here, we focus on text summarization and exclude multi-modal summarization systems, such as source code (Ahmad et al., 2020), and screen-play summarization (Papalampidi et al., 2019, 2020).

First, as expected, ROUGE is used by more than 95% of papers, while other metrics such as METEOR, BLEU, BERTSCORE, and MOVERSCORE are rarely used. Interestingly, 64% of the surveyed papers used manual evaluation to analyze the strengths and weaknesses of the proposed model(s), an area where the single figure-of-merit output of ROUGE does not provide direct insights.

In Table 1, we summarize the 6 major dimensions of manual evaluation as faithfulness, recall, precision, relevance, coherence, and fluency. Faithfulness is the degree of factual consistency with respect to the source article. Recall, precision, and relevance measure the degree of salient and important information, where relevance is generally measured as the harmonic mean of precision and recall. According to our analysis, recall, precision, faithfulness, and fluency are the most frequent dimensions of human evaluation in recent work, from which we take inspiration in designing FFCI (with fluency defined as inter-sentential coherence, as discussed in Section 1).

We also found that absolute scoring is more common than relative evaluation. Absolute benchmark is conducted by asking annotators to evaluate system-generated summaries based on a numeric scale, in isolation of any other summaries. With relative evaluation, on the other hand, annotators are asked to directly rank summaries generated by different methods.

Lastly, the manual evaluations in Table 1 were conducted with different basis, namely, SCU (semantic content units, as defined in Pyramid, Nenkova & Passonneau, 2004), reference, article, and reference+article. SCU is clauses or sentences that are manually extracted from the ground-truth summary, and are used to evaluate content selection in summarization. Pyramid method is initially applied to aggregate the human summaries, however, previous work (Bhandari et al., 2020) applied Pyramid in the single-reference setting.

2. The first 26 rows are from Hardy et al. (2019). We manually re-examine these papers and found miss-annotation for Amplayo et al. (2018).

We observe that most recent work has used the source article as the basis in assessing faithfulness, precision, and recall, rather than reference summaries or SCUs. Intuitively, this is the best practice in human evaluation, especially for faithfulness, as generated summaries can technically contain details not found in reference summaries but are in the source article, and they should still be seen as faithful information in this case. However, for precision and recall, Nenkova and Passonneau (2004), Fabbri et al. (2020) have shown that using the source article rather than reference summaries leads to poor inter-annotator agreement as a result of the complication in the annotation scheme. Although most papers of the 71 papers in Table 1 that perform manual evaluation base the evaluation on the source article, only 7 out of 71 papers describe explicit quality control mechanisms used in their experiments.³

2.2 Resource for Summarization Evaluation

Best practice in assessing quality automatic metrics for text generation systems is by measuring correlation scores such as Pearson, Spearman, or Kendall between system-generated text and a reference. In machine translation (MT), BLEU (Papineni et al., 2002) and METEOR (Lavie & Agarwal, 2007) were validated based on WMT and LDC TIDES 2003 corpora, respectively. While MT metric evaluation resources have been developed progressively over time (e.g. the WMT Metrics Task has run annually since 2006), there has been a relative dearth of new evaluation datasets for summarization research, and only recently have Bhandari et al. (2020) and Fabbri et al. (2020) released evaluation datasets based on summaries generated by a range of neural summarization models.

Table 2 comprehensively lists the available resources for summarization evaluation research, in which we observe an 11 year gap between DUC-TAC and the recent datasets. Because the summaries in the DUC⁴ and TAC⁵ datasets are from more than 10 years ago, they are based on largely outdated extractive summarization systems. Bhandari et al. (2020), Fabbri et al. (2020), Maynez et al. (2020), Wang et al. (2020) attempt to tackle this issue by releasing new data, although the dimensions of evaluation represented in those datasets do not fully align with the common dimensions of manual evaluation in Table 1. For instance, Bhandari et al. (2020) only assess coverage based on SCUs, and Fabbri et al. (2020) do not separate out precision and coverage.

Bhandari et al. (2020) annotated 100 samples based on the simplified Pyramid method (Nenkova & Passonneau, 2004), where semantic content units (SCUs) are manually extracted and crowd-workers then count the appearance of SCUs in the summary. This annotation scheme is closely related to coverage as proposed in this research, but does not consider focus, faithfulness, and inter-sentential coherence. Bhandari et al. (2020) and Peyrard (2019b) both found that evaluation metrics developed based on older datasets do not necessarily perform well on modern datasets with more modern summarization systems.

Fabbri et al. (2020) assess four dimensions of summaries: relevance, consistency, fluency, and coherence, by annotating 100 CNN/DailyMail samples. Our FFCI framework further decomposes relevance into focus and coverage to provide a more fine-grained understanding of content overlap, and replaces fluency — which measures quality of individual sentences

3. Quality control is a mechanism to measure the quality of the crowd-sourced annotation (Graham et al., 2016).

4. <https://duc.nist.gov/data.html>

5. <https://tac.nist.gov/data/>

Dataset		#Systems	#Summaries per System	Aspects	Annotation benchmark
DUC-2001	SDS	12	149	Coverage	1 reference
DUC-2002		14	295		
DUC-2003		14	624		
DUC-2001	MDS	42	29	Coverage	1 reference
DUC-2002		36	59		
DUC-2003		18	30		
TAC-2008	SDS	54 + 4 ref	480	Pyramid (Coverage) Responsiveness	1 reference
TAC-2009		55 + 4 ref	440		
Bhandari et al. (2020) (CNNDM)	SDS	25	100	Pyramid (Coverage)	1 reference
Fabbri et al. (2020) (CNNDM)		23	100		
Maynez et al. (2020) (XSUM)	SDS	4	500	Faithfulness	article
Durmus et al. (2020) (CNNDM, XSUM)		N/A	1034*		
Wang et al. (2020) (CNNDM, XSUM)		1	474		

Table 2: Resources for summarization evaluation. MDS/SDS = Multi/Single Document Summarization. * indicates the total of summaries for all system as #Systems is not reported by Durmus et al. (2020).

— with inter-sentential coherence, which measures the quality of multi-sentence summaries more holistically. They evaluated summaries via crowd-sourcing (Amazon MTurk) and expert (in-house) annotators, but ultimately base all of their findings on the expert annotations, as they found the crowd-sourced annotations to be highly inconsistent. First, their annotation scheme is difficult for crowd-workers, as they are asked to judge all four dimensions after reading an article and a system-generated summary. Consistency (faithfulness) is found to be particularly difficult (and subjective), and previous studies (Maynez et al., 2020) have attempted to ease the annotation burden by asking crowd-workers to highlight unfaithful spans in the summary. Assessing relevance without a ground-truth summary is also hard, as it requires crowd-workers to implicitly construct their own summary of the article. The second reason is that there is no quality control to verify the quality of the annotations, which means they may be potentially unreliable. In this work, we use the resource released by Maynez et al. (2020) to study faithfulness, and use the customized *Direct Assessment* framework (Graham et al., 2015) to collect judgements across three additional dimensions: focus, coverage, and inter-sentential coherence. The annotation framework we use has the following benefits: a more intuitive annotation scheme, better quality control, and better handling of annotator variance (through z -score normalization).

Perhaps more importantly, despite presenting extensive evaluation on a number of state-of-the-art summarization systems using a wide range of evaluation methods, Fabbri et al. (2020) stop short of providing guidance as to the best evaluation methods for assessing a particular dimension of summary quality. Our work addresses this gap by providing practical advice on the best evaluation method for assessing the four dimensions of FFCI.

3. Existing Evaluation Metrics and Extensions

In this section, we review the existing evaluation metrics that we use for different dimensions of summarization evaluation (faithfulness, focus, coverage, and inter-sentential coherence). We first introduce traditional string overlap-based evaluation metrics, with a particular focus on ROUGE (Lin, 2004), as it has become the de facto standard for automatic summarization evaluation. We apply the overlap-based metrics in evaluating all four dimensions of FFCI (see Table 1). Next, we present the recently-proposed QAGS question answering-based framework for evaluating faithfulness (Wang et al., 2020), which we extend to also evaluate focus and coverage (but not inter-sentential coherence). We then introduce two general-purpose string similarity metrics, namely the unsupervised BERTSCORE (Zhang et al., 2020b) and supervised STS-SCORE, which is trained over STS data from successive SemEval tasks (Agirre et al., 2012). Both of these metrics are used to evaluate all four dimensions of FFCI. Finally, we introduce the coherence score of Nayeem and Chali (2017) and Yin et al. (2020) as a specialized metric for evaluating inter-sentential coherence.

3.1 Traditional String Overlap-Based Evaluation Metrics

Despite its brittleness, the simplicity of ROUGE (Lin, 2004) has made it a mainstay of summarization evaluation for over 15 years. ROUGE measures the overlap between a generated and reference summary in terms of unigram or bigram overlap (ROUGE-1 and ROUGE-2, respectively), or longest common subsequence (ROUGE-L). In another work, Ng and Abrecht (2015) proposed ROUGE-WE as an extension of ROUGE which incorporates word2vec (Mikolov et al., 2013) embeddings, but found it to perform similarly to the simpler ROUGE-1 and ROUGE-2 metrics in practice.

In most studies, the harmonic mean (F1) of each of the three main ROUGE variants (ROUGE-1, ROUGE-2 and ROUGE-L) is used for evaluation, which some studies (See et al., 2017; Pasunuru & Bansal, 2018) complement with machine translation metrics such as BLEU (Papineni et al., 2002) and METEOR (Lavie & Agarwal, 2007).

3.2 Question Answering-Based Evaluation

Recent work (Wang et al., 2020; Durmus et al., 2020) has shown that factual consistency can be evaluated using a question answering (“QA”) task formulation. In this paper, we experiment with the QAGS framework (Wang et al., 2020), which involves two components: (1) question generation (“QG”), and (2) question answering.

Let X , Y , and Y' be the source document, reference summary, and system summary, respectively. For faithfulness, QAGS defines $p(Q|Y')$ as the distribution over questions Q generated from system summary Y' . Answer A is predicted based on two terms: $p(A|Q, X)$ and $p(A|Q, Y')$, representing the answer distribution based on the source document and system summary, respectively. Factual consistency is measured by the F1 score (or exact match) between the answers generated from the source document and system summary.

Evaluating faithfulness via QA-based evaluation such as QAGS is intuitive but has several drawbacks. First, QAGS requires careful tuning of hyperparameters such as the number of questions to generate, maximum token length for question generation, and also question filtering method. Secondly, QA-based evaluation is computationally expensive and

hard to apply to languages other than English, due to the need for training data for the QA and QG models.

Despite its drawbacks, previous work has reported encouraging results in the evaluation of faithfulness. In this work, we extend the QA-based method to two other elements of FFCI: focus and coverage. In this, we address the following research question: **[RQ1]** *how effective is QAGS relative to other simpler methods for assessing faithfulness, and can it be applied to evaluate focus and coverage?*

3.3 BERTScore

Contextualized word embeddings have been shown to be a strong metric for evaluating machine translation (Mathur et al., 2019; Zhang et al., 2020b). Zhang et al. (2020b) proposed BERTSCORE as a means of computing the similarity between BERT token embeddings of system and reference texts, while in other work Zhao et al. (2019) proposed MOVERSCORE as the Euclidean distance between two contextualized BERT representations. We use BERTSCORE rather than MOVERSCORE in this study for two reasons: (1) MOVERSCORE is symmetric (i.e. $\text{MoverScore}(x, y) = \text{MoverScore}(y, x)$), and as such cannot easily be used to evaluate precision and recall separately; and (2) recent work (Fabbri et al., 2020) has shown that BERTSCORE is superior to MOVERSCORE for summarization evaluation.

For Y and Y' as the reference and system summary, in the context of summarization, BERTSCORE is computed as follows:

$$\begin{aligned}\mathcal{P}_{\text{BERT}} &= \frac{1}{|Y'|} \sum_{t_i \in Y'} \max_{s_j \in Y} t_i^T s_j \\ \mathcal{R}_{\text{BERT}} &= \frac{1}{|Y|} \sum_{s_j \in Y} \max_{t_i \in Y'} t_i^T s_j \\ \mathcal{F}_{\text{BERT}} &= 2 \frac{\mathcal{P}_{\text{BERT}} \cdot \mathcal{R}_{\text{BERT}}}{\mathcal{P}_{\text{BERT}} + \mathcal{R}_{\text{BERT}}}\end{aligned}$$

where s_j and t_i are token embeddings of Y and Y' .

In terms of hyperparameters, BERTSCORE is simpler than QA-based evaluation, with the main hyperparameter being layer selection: Zhang et al. (2020b) found that selection of which transformer layer to source the token embeddings from is critical to performance. For machine translation and text generation evaluation, Zhang et al. (2020b) recommend the use of $\mathcal{F}_{\text{BERT}}$ based on the 24th layer of `roberta-large`, on the basis of experiments over BERT (Devlin et al., 2019), RoBERTA (Liu et al., 2019b), and XLNET (Yang et al., 2019).

Since layer selection in the original paper was based on machine translation datasets, we perform similar layer selection across the three sub-facets of FFCI, asking: **[RQ2]** *which layer of which pre-trained language model is best for evaluating faithfulness, focus, and coverage of a summary?*

We perform a model-layer search to answer this question, extending the work of Zhang et al. (2020b) to include other pre-trained models. In total, we examine 7 model types, that can be categorized as follows: (1) encoder-only = BERT (Devlin et al., 2019), RoBERTA (Liu et al., 2019b), and XLNET (Yang et al., 2019); (2) decoder-only = GPT2 (Radford

et al., 2019); and (3) encoder–decoder = T5 (Raffel et al., 2019), BART (Lewis et al., 2020), and PEGASUS (Zhang et al., 2020a). For each of these, we experiment with different-sized pre-trained models, for a total of 19 models. For encoder–decoder models, we only perform layer selection over the encoder layers.

3.4 Model-Based Approaches

3.4.1 FACTCC

Kryscinski et al. (2020) proposed a weakly-supervised, model-based approach for verifying factuality in abstractive summaries. The training data is generated based on transformation rules including paraphrasing, entity and number swapping, pronoun swapping, sentence negation, and noise injection. The goal is to estimate $P(y|A, c)$ where y is a binary label of **CORRECT** and **INCORRECT**, A is the source article, and c is the transformed sentence (claim or summary). For training, Kryscinski et al. (2020) simply fine-tuned BERT model and use [CLS] for classification (and denote this model as FactCC). Additionally, the training is extended by allowing the model to not only classify the claim consistency but also highlight a span in the source article as the supporting evidence (and denote this model as FactCCX).

3.4.2 STS-SCORE

We additionally experiment with STS-SCORE. Semantic textual similarity (STS) measures the relative semantic similarity of two short texts (often single sentences) on a continuous scale of $[0, 5]$ (Agirre et al., 2012). A broad range of STS approaches have been proposed, and datasets have been released for a number of different languages, predominantly through successive SemEval tasks.

Similar to BERTSCORE, STS-score is a similarity function, from which precision, recall, and F1 can be calculated as follows:

$$\begin{aligned} \mathcal{P}_{\text{STS}} &= \frac{1}{|Y'|} \sum_{t_i \in Y'} \max_{s_j \in Y} \text{STS}(t_i, s_j) \\ \mathcal{R}_{\text{STS}} &= \frac{1}{|Y|} \sum_{s_j \in Y} \max_{t_i \in Y'} \text{STS}(s_j, t_i) \\ \mathcal{F}_{\text{STS}} &= 2 \frac{\mathcal{P}_{\text{STS}} \cdot \mathcal{R}_{\text{STS}}}{\mathcal{P}_{\text{STS}} + \mathcal{R}_{\text{STS}}} \end{aligned}$$

where s_j and t_i are segments within reference summary Y and system summary Y' , respectively. We experiment with three segment granularities: (1) elementary discourse units (EDUs) (Mann, 1984);⁶ (2) sentence; and (3) document. As our STS scorer, we use a fine-tuned sentence transformer, based on the findings of Reimers and Gurevych (2019a).

3.5 Coherence Score

Nayeem and Chali (2017) and Yin et al. (2020) define coherence score as the weighted sum of similarity scores of two adjacent sentences. For system summary Y' , the coherence score

6. Typically a clause or sentence that represents atomic information in discourse parsing. Reference and generated summaries may have different sentence lengths/granularities, and EDU-based segmentation is a possible alternative to sentence-based matching

Context	Maynez et al. (2020)	Durmus et al. (2020)	Wang et al. (2020)
XSUM samples	2000 (4 models)	286 (unknown)	239 (1 model)
CNN/DailyMail samples	—	748 (4 models)	235 (1 model)
Filtered data sampling	No	Yes (meaningful sentence)	No
Quality control in annotation	Yes (pilot study)	None	Yes
Evaluation against reference	ROUGE, BERTSCORE	ROUGE	ROUGE, BERTSCORE, BLEU, METEOR
Evaluation against article	QA, Entailment	ROUGE, BLEU, BERTSCORE, Entailment, QA	QA
Best evaluation	Entailment	QA	QA

Table 3: Summary of previous work on faithfulness evaluation.

is computed as follows:

$$\text{coherence}(Y') = \frac{1}{n-1} \sum_{i=1}^{n-1} \text{Sim}(t_i, t_{i+1})$$

$$\text{Sim}(t_i, t_{i+1}) = \lambda \text{NESim}(t_i, t_{i+1}) + (1 - \lambda) \text{CosSim}(t_i, t_{i+1})$$

where NESim is the named entity overlap of two sentences t_i and t_{i+1} in Y' , and cosine similarity is measured based on pre-trained word embeddings.⁷ While this method is commonly used to assess coherence in the sentence ordering task (Shen & Baldwin, 2021), this is the first paper to systematically evaluate its effectiveness for summarization evaluation.

4. FFCI Framework

4.1 Faithfulness

Abstractive summarization is prone to “hallucination” or factual inconsistencies, where information is generated that does not exist in the source document (Maynez et al., 2020; Wang et al., 2020). Three recent papers independently proposed to evaluate the degree of hallucination (Maynez et al., 2020; Durmus et al., 2020; Wang et al., 2020), as detailed in Table 3.

In terms of training data, Maynez et al. (2020) released the largest dataset with 2,000 annotated summaries generated over XSUM. Durmus et al. (2020) manually pre-filtered the data to select “meaningful” sentences, making it difficult to fully automate the method, and do not report on any quality controls in their human annotation. On the other hand, Maynez et al. (2020) conducted a pilot study to train their annotators, and Wang et al. (2020) applied annotator attention checks, making us more confident in the quality of the resultant dataset.

7. We use GloVe embeddings (Pennington et al., 2014) to compute coherence scores, consistent with Nayeem and Chali (2017).

Metric	Question	Answer
Faithfulness	$p(Q Y')$	$p(A Q, X)$ and $p(A Q, Y')$
Focus	$p(Q Y')$	$p(A Q, Y)$ and $p(A Q, Y')$
Coverage	$p(Q Y)$	$p(A Q, Y)$ and $p(A Q, Y')$

Table 4: Probability distributions used in QAGS to evaluate faithfulness, focus, and coverage.

We argue that the best way to evaluate faithfulness is by comparing the generated summary with the source document (and not with the reference summary).⁸ In Table 3, Durmus et al. (2020) is the only paper to extensively measure traditional and model-based metrics against the source article. However, because of concerns over their data, we revisit faithfulness evaluation using the dataset of Maynez et al. (2020).⁹ We score faithfulness by comparing summary sentences and the source document as follows:

$$\text{FA}_{\text{METRIC}} = \frac{1}{|Y'|} \sum_{t_i \in Y'} A(t_i, X, n)$$

$$A(t_i, X, n) = \text{AvgTop-}n \text{ METRIC}(t_i, s_j)$$

where t_i and s_j are sentences from the system summary Y' and source document X , respectively; $\text{METRIC} \in \{\text{ROUGE}, \text{STS-SCORE}, \text{BERTSCORE}\}$; and $n \in \mathbb{Z}^+$ is a hyperparameter. $\text{AvgTop-}n$ matches sentence t_i from the summary with each sentence s_j in the source document X , and returns the average score for the top- n best-matching sentences. The intuition behind measuring across the top- n is that information in a summary sentence might potentially be drawn from different sentences in the source article.

For faithfulness, ROUGE, STS-SCORE, and BERTSCORE are based on F1-scores. In preliminary experiments, we compared $n \in \{1, 2, 3\}$ and found that $n = 2$ works best for ROUGE, and $n = 3$ works best for STS-SCORE and the pre-trained language model scores.

4.2 Focus and Coverage

We test the ability of the evaluation metrics from Section 3 to measure focus and coverage, arguing that it is important to separately measure summaries in terms of precision and recall relative to a reference.

First, we adopt QAGS from faithfulness evaluation (Wang et al., 2020), and extend it to evaluate focus and coverage based on the probability distributions in Table 4. For focus, we generate questions Q based on system summary ($p(Q|Y')$) similarly to faithfulness, and answer the questions based on $p(A|Q, Y)$ and $p(A|Q, Y')$. That is, we test the consistency of answers generated from the system and reference summaries, based on questions generated from the system summary (meaning we only evaluate information present in the *system*

8. As we argued earlier (Section 2.1), details in the generated summary that are not in the reference summary but in the source article should still be regarded as faithful information.

9. At the time this research was conducted, only Maynez et al. (2020) had released their data.

summary as this is the source of the questions; hence focus). For coverage, on the other hand, we generate questions based on the *reference* summary ($p(Q|Y)$), and answer those questions based on $p(A|Q, Y)$ and $p(A|Q, Y')$, in the same manner as focus (meaning we evaluate information present in the *reference* summary; hence coverage). We return to discuss how to generate questions and answers in Section 5.3.

Apart from QAGS, we also examine ROUGE, METEOR, BLEU, BERTSCORE, and STS-SCORE to evaluate focus and coverage. For computing ROUGE, STS-SCORE, and BERTSCORE, we use the precision and recall for focus and coverage, respectively.

4.3 Inter-Sentential Coherence

We extend the Nayeem and Chali (2017) method to measure inter-sentential coherence (IC) within system summary Y' , based on a next-sentence-prediction (NSP) classifier as follows:

$$\text{NSP}(Y') = \text{mean}_{t_i \in Y'} \text{NSP}(t_i, t_{i+1})$$

where each t_i is a sentence in summary Y' , and NSP returns a probability of t_{i+1} following t_i . We experimented with max, min, and mean aggregation, but found mean to produce the most robust results.

Compared to Nayeem and Chali (2017), our NSP score also assesses coherence between two adjacent sentences, but with a model-based system, rather than based on cosine similarity of pre-trained word embeddings. This way, we argue that the NSP score can better assess the overall writing flow based on two adjacent sentences because it is not limited by the factual content of the sentences.¹⁰

Note that we do not use the NSP classifier in pretrained language models, as not all pretrained language models have this objective. Instead, we train a separate NSP classifier (by fine-tuning pretrained language models) where positive examples are two consecutive sentences and negative samples are constructed using a range of strategies (e.g. by flipping the sentences), as detailed in Section 5.3.

In contrast to ROUGE, STS-SCORE, and BERTSCORE, our proposed evaluation scheme for inter-sentential coherence is a reference-less metric. ROUGE, STS-SCORE, and BERTSCORE are reference-based metrics and designed for evaluating saliency and coverage when compared to a reference text. That said, reference-based metrics might implicitly assess inter-sentential coherence because a system summary that is similar to the gold summary is likely to be coherent (since the human-written gold summary should be coherent).

5. Experimental Setup

5.1 Data

In order to evaluate the different metrics and perform model-layer selection, we need gold-standard data for each of the four FFCI sub-tasks.

10. Having said which, although the facts in both sentences can be different, they should still have similar topics and flow coherently.

FAITHFULNESS

We use the 2000 samples from Maynez et al. (2020), which is based on summaries generated by 4 neural models over XSUM (Narayan et al., 2018b): pointer generator network (“PG”: See et al. (2017)), Topic-aware convolutional Seq2Seq (“TCNV”: Narayan et al. (2018b)), a transformer-based model (“TRANS2S”: Vaswani et al. (2017)), and BERT (“BERT”: Devlin et al. (2019), Liu and Lapata (2019)).¹¹

FOCUS AND COVERAGE

We annotate 1080 data-model pairs by randomly sampling 135 articles each from the test sets of CNN/DailyMail (Hermann et al., 2015) and XSUM (Narayan et al., 2018b), and generate summaries with two models: PG (See et al., 2017) and BERT (Liu & Lapata, 2019); this results in 540 summaries ($135 \times 2 \times 2$) which are assessed for focus and coverage.

Note that for CNN/DailyMail, we use the PG+Coverage variant, while for XSUM, the basic PG model is used, as it produces better summaries (See et al., 2017; Narayan et al., 2018b). We choose these data-model pairs for two reasons: (1) CNN/DailyMail and XSUM are benchmark abstractive summarization datasets which represent the most extractive and abstractive summarization corpora, respectively (Bommasani & Cardie, 2020)¹²; and (2) PG and BERT are representative of contemporary neural models from the attention-based recurrent model to pre-trained language model era.

INTER-SENTENTIAL COHERENCE

We used the same 270 system summaries from CNN/DailyMail as for focus and coverage.¹³

5.2 Human Evaluation Task Design

We used Amazon Mechanical Turk¹⁴ and the customized Direct Assessment (“DA”) method (Graham et al., 2015; Graham et al., 2017), which has become the de facto for MT evaluation in WMT. DA equips the annotation scheme with some pre-annotated samples for quality control, two texts (system and human translation), and a slider button (continuous scale with range 1–100) for annotation. In Figure 2, we present the annotation interface of the customized DA method for summarization evaluation.

For focus and coverage, the annotation interface provides system and reference summary, and a question: *How much information contained in the second text can also be found in the first text?* We were able to combine focus and coverage annotation, as the only thing that differentiates them is the ordering of the system and reference summaries, which was invisible to annotators.¹⁵ For inter-sentential coherence, the annotators are given a single summary and asked to rate inter-sentential coherence directly.

11. Note that, at the time of writing, there is no annotated dataset for faithfulness based on CNN/DailyMail, so we can only evaluate faithfulness over XSUM.

12. Noting that contemporaneous works (Bhandari et al., 2020; Fabbri et al., 2020) only use CNN/DailyMail.

13. Noting that XSUM summaries are single sentences, and thus inter-sentential coherence is not relevant.

14. <https://www.mturk.com/>

15. For focus, the first and second texts are the reference and system summaries, respectively. For coverage, the order is reversed, and they are the system and reference summaries, respectively.

How much information contained in the black text can also be found in the gray text?

officials at the famous yellowstone national park in the us have revealed that they had to put down a newborn bison after some tourists put it in the boot of their car .

wildlife rangers in the us state of wyoming have warned visitors to stay away from their herd after they refused a controversial bison .

0 % 100 %

(a) Focus and coverage annotation.

The text below has good inter-sentential coherence (i.e. the flow from one sentence to the next is natural):

sony emails reveal bbc bosses want to turn the hit series starring peter capaldi and is screened in 50 countries , to be turned into a movie to capitalise on its worldwide success .

but the emails show doctor who's creative team are reluctant to rush into making a film that could flop and tarnish its reputation .

Strongly disagree Strongly agree

(b) Inter-sentential coherence annotation.

Figure 2: MTurk annotation interface for focus, coverage, and inter-sentential coherence.

While it may seem more natural and reliable to evaluate focus and coverage based on the source document than the ground-truth summary, we use the ground-truth summary in this research for the following reasons. First, historically, validation of automatic summarization evaluation metrics has been based primarily on ground-truth summaries (not source documents). Second, previous work such as DUC (dataset for ROUGE), TAC (dataset for MOVERSCORE), and Bhandari et al. (2020) annotated coverage based on a single reference summary. Third, this work is based on single-document summarization systems, and we argue that the variance in content is actually not that great. Lastly, basing human evaluation (of focus and coverage) on the source article leads to more complicated annotation schemes, and has been shown to yield poor annotations (as discussed in Section 2.1).

We posted separate HITs for focus + coverage vs. inter-sentential coherence, where a single HIT consisted of 100 annotation instances including 10 quality control instances. For focus + coverage, 5 samples are random pairs (should be scored 0) and the remaining

Aspect	Pearson correlation (r) (agreement)	Avg. Quality score (%)	Avg. Working time (min)
Focus + Coverage	0.57	90.3	62.1
Inter-sentential coherence	0.49	94.4	35.8

Table 5: Statistics over the approved HITs on both criteria: Focus + Coverage and Inter-sentential coherence. Quality score is the average score of 10 quality control samples for each HIT.

5 samples are repetitions with minor edits (should be scored 100). For inter-sentential coherence, 5 samples are random sentence pairs, and the remaining 5 are verbatim repeats (both of which should be scored 0).

We restricted the HITs to US-based workers with at least 10,000 approved HITs. For each HIT, we pay USD\$5 and an additional \$8 bonus if they pass the quality control checks (the minimum bar is to have at least 7 correct answers from the total of 10 quality control samples), to ensure that workers are paid at a level that is comfortably above the minimum wage in Australia.¹⁶

We collected 3 annotations per HIT which passed quality control (running new HITs in cases where HITs had to be discarded), and present the statistics over HITs in Table 5. We achieved a mean Pearson’s correlation between annotators of $r = 0.57$ and 0.49 , for focus + coverage and inter-sentential coherence, respectively. We also observe that the average quality score is high and the working time of both HITs is reasonable. To aggregate the scores, we standardized the scores of each worker into a z -score before averaging.

5.3 Evaluation Metrics

5.3.1 ROUGE, METEOR, AND BLEU

In this experiment, we use the original implementation of ROUGE¹⁷ and METEOR¹⁸. For BLEU, it is based on SacreBLEU implementation (Post, 2018).¹⁹

5.3.2 QAGS

We re-implemented QAGS (Wang et al., 2020) by training the question generator with `bart-large` on NewsQA (Trischler et al., 2017), and QA model with `bert-large-wwm` on SQuAD2.0 (Jia et al., 2018), achieving similar results to Wang et al. (2020) on both tasks. We generate a maximum of 50 questions, and discard questions if the QA system cannot predict the correct answer based on the original context the question is generated from.

16. For inter-sentential coherence annotation, we pay USD\$3 and an additional \$3 bonus.

17. <https://github.com/bheinzerling/pyrouge>

18. <http://www.cs.cmu.edu/alavie/METEOR/>

19. <https://github.com/mjpost/sacrebleu>

To validate our implementation, we tested the model at faithfulness evaluation over XSUM using the dataset of Maynez et al. (2020), and achieved a correlation of $r = 0.25$, 0.075 points higher than the original paper.²⁰

5.3.3 BERTSCORE

We sourced 19 pre-trained language models from HuggingFace,²¹ and adjust the BERTSCORE implementation to search for the best model-layer combination.²² The models include BERT (Devlin et al., 2019), ROBERTA (Liu et al., 2019b), XLNET (Yang et al., 2019), GPT2 (Radford et al., 2019), T5 (Raffel et al., 2019), BART (Lewis et al., 2020), and PEGASUS (Zhang et al., 2020a).²³

In selecting the best layers of each model for faithfulness, focus, and coverage, we do not merge datasets and systems, but instead select based on the averaged best results across different dataset-system combinations. We believe this is a robust method which does away with the need for a method such as cross-validation. Another reason not to merge datasets and systems (and compute correlation across summary-level data points) is because of the small number of systems: unlike DUC and TAC which have many systems, the faithfulness data (Maynez et al., 2020) only consists of 4 different systems from 1 dataset, while our focus and coverage data consist of 2 systems from 2 datasets.

5.3.4 FACTCC

We use the models and original implementation of FactCC and FactCCX from Kryscinski et al. (2020) to evaluate faithfulness.²⁴ We use the probability of class CORRECT (a summary being factually correct relative to the source article) as the final output for both FactCC and FactCCX.

5.3.5 STS-SCORE

We use sentence-transformers (Reimers & Gurevych, 2019b) with `bert-large-nli`, using `spacy`²⁵ and discourse segmentation (Ji & Eisenstein, 2014) to perform sentence and EDU segmentation. We also experimented with other pre-trained transformer models, but found there to be little difference in the results.

5.3.6 NSP SCORE

For inter-sentential coherence, we fine-tune a pretrained language model for NSP classification on 100,000 sentence pairs (50K positive and 50K negative) automatically extracted from original XSUM articles. We experimented with four types of negative samples: `type1`

20. Noting that at the time of writing, the QAGS data and code have not been released, and that our evaluation is thus based on a different test dataset to the authors.

21. <https://huggingface.co/>

22. https://github.com/Tiiiger/bert_score

23. `bert-base-uncased`, `bert-large-uncased`, `roberta-base`, `roberta-large`, `roberta-large-mnli`, `xlnet-base-cased`, `xlnet-large-cased`, `gpt2`, `gpt2-medium`, `gpt2-large`, `gpt2-xl`, `t5-small`, `t5-base`, `t5-large`, `bart-base`, `bart-large`, `pegasus-xsum`, `pegasus-cnn_dailymail`, `pegasus-large`.

24. <https://github.com/salesforce/factCC>

25. <https://spacy.io/>

Data Variant	C-PG	C-BT
1	0.22±0.06	0.27±0.01
2	0.23±0.02	0.28±0.04
3	0.25±0.12	0.30±0.12
4	0.41±0.01	0.26±0.07
5	0.39±0.07	0.35±0.05

Table 6: Averaged Pearson correlation scores for inter-sentential coherence and NSP-Score over 5 run models for the C-PG (CNN/DailyMail-PG) and C-BT (CNN/DailyMail-BERT) data, based on the five training data variants.

= flipped sentence pairs; **type2** = pairs where the second sentence is randomly obtained from a different document; **type3** = pairs of corrupted repetitive sentences; and **type4** = pairs where the second sentence is randomly picked from the same document in arbitrary position.

We define 5 training data variants by combining the types as follow: (1) $50\text{K} \times \text{type1}$; (2) $50\text{K} \times \text{type2}$; (3) $25\text{K} \times \text{type1} + 25\text{K} \times \text{type2}$; (4) $25\text{K} \times \text{type2} + 5\text{K} \times \text{type3} + 20\text{K} \times \text{type4}$; and (5) $25\text{K} \times \text{type1} + 5\text{K} \times \text{type3} + 20\text{K} \times \text{type4}$.

In preliminary experiments, we tested seven models (BERT, RoBERTa, ALBERT, XLNet, ELECTRA, GPT2, and BART) for fine-tuning NSP score. However, the results indicated that BERT performs the best for NSP score (see Appendix D), so this forms the basis of our primary results in the paper for inter-sentential coherence. First, we partition our data into training, development, and test splits based on a ratio of 80:10:10, respectively, and fine-tune `bert-base-uncased` with learning rate = $5e-5$, batch size = 40, and maximum epochs = 20. We simply use the [CLS] encoding as the input to an MLP layer. During training, we use early stopping (patience = 5) based on the development set performance. We run all models 5 times and achieve varied averaged F1 scores ranging from 75% to 93%, but found variant-5 to achieve the best overall Pearson correlation (see Table 6).

6. Experimental Result

6.1 Faithfulness

In Table 7, we show Pearson correlation scores (r) for faithfulness in two different forms: (1) evaluation against the reference; and (2) evaluation against the source. We additionally report the Spearman correlation (ρ), for direct comparison with Maynez et al. (2020). We measure the correlation between human judgement (data in Section 5.1) and various automatic metrics. For evaluation against the reference, FA_{ROUGE} and $\text{FA}_{\text{BERTSCORE}}$ are equivalent to ROUGE and BERTSCORE, respectively, because the XSUM dataset only consists of one-sentence summaries.

Metric	r	ρ
<i>Against reference</i>		
ROUGE-1	0.199	0.199
ROUGE-2	0.116	0.161
BLEU-4	0.072	0.133
METEOR	0.131	0.170
BERTSCORE*	0.128	0.131
<i>Against source sentences</i>		
ROUGE-1	-0.047	-0.028
ROUGE-2	0.179	0.221
BLEU-4	-0.107	-0.138
METEOR	-0.018	0.006
BERTSCORE*	-0.034	0.006
FactCC	0.042	0.045
FactCCX	-0.027	-0.017
QA (Maynez et al., 2020)	—	0.044
Entailment (Maynez et al., 2020)	—	0.431
QAGS (our implementation)	0.250	0.270
FA _{STS}	0.260	0.258
FA _{ROUGE-1}	0.361	0.361
FA _{ROUGE-2}	0.311	0.315
FA _{BERTSCORE} *	0.178	0.179
FA _{BERTSCORE} (OURS)	0.476	0.474

Table 7: Pearson (r) and Spearman (ρ) correlation coefficients for faithfulness, measured between human judgement and various automatic metrics. “*” denotes that BERTSCORE uses `roberta-large` (layer 24 as recommended by Zhang et al. (2020b)) while ours uses `roberta-base` (layer 10).

As we can see in Table 7, computing ROUGE, BLEU, METEOR and BERTSCORE conventionally (either against reference or full source sentences) yields low correlation scores.²⁶ We also found that our QAGS implementation (Wang et al., 2020) performs better than the QA system of Maynez et al. (2020) over their dataset. For FactCC and FactCCX, they have low correlation scores, with 0.04 being the highest correlation score. This is in line with a recent study by Pagnoni et al. (2021) who found that FactCC performs poorly (0.07 correlation score) over the XSUM dataset.

When we apply ROUGE, STS-SCORE, and BERTSCORE over the source document based on FA_{METRIC} in Section 4.1, we see that the FA_{ROUGE-1} and FA_{ROUGE-2} baselines actually outperform QAGS that is computationally expensive, but more importantly that the two versions of FA_{BERTSCORE} perform differently, with our summarization-optimized version resulting in the best overall results. The first FA_{BERTSCORE} uses the recommendations of

26. For evaluation against the reference, we reproduce the Spearman correlation scores (ρ) of ROUGE-1 and ROUGE-2, but BERTSCORE is 0.131, slightly lower than Maynez et al. (2020) ($\rho = 0.190$). We use the recommended model-layer by Zhang et al. (2020b) while Maynez et al. (2020) do not report their BERTSCORE configuration (e.g. source code or model-layer used for evaluation).

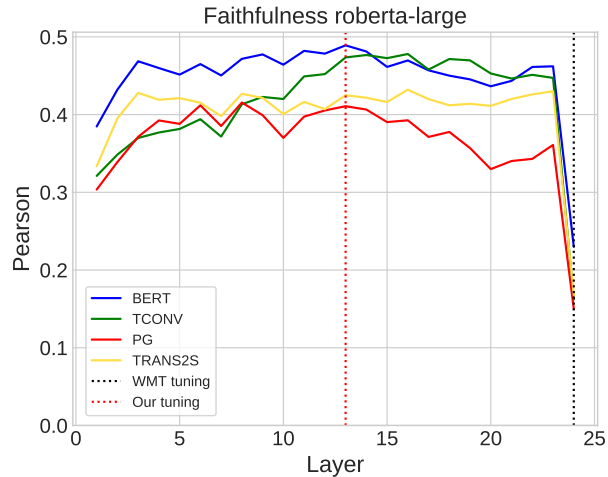


Figure 3: Pearson correlation based on each layer of `roberta-large` for faithfulness evaluation. Tuning refers to layer selection (i.e. model parameters are not updated.)

Zhang et al. (2020b) (`roberta-large`; layer 24), but we found this configuration to produce a lower correlation. The best layer for this model is layer-13, as depicted in Figure 3.

After conducting model-layer search over 19 models,²⁷ we found that `roberta-base` (layer 10) to result in the best correlation (based on average rank across four summarization models).

Matching information in the summary sentence with the article sentence works best through $FA_{\text{BERTSCORE}}$ in Table 7. Although $FA_{\text{BERTSCORE}}$ is computationally cheaper than question-answering based models such as QA and QAGS, we argue that deeper investigation is needed to more thoroughly evaluate $FA_{\text{BERTSCORE}}$, especially on different dataset with more varied article and summary lengths. As this paper is focused on introducing the FFCI framework, for faithfulness we use only the available data from Maynez et al. (2020), and leave further validation to future work.

6.2 Focus, Coverage, and Inter-Sentential Coherence

6.2.1 DATASET VISUALIZATION

In Figure 4, we visualize the focus and coverage data for CNN/DailyMail and XSUM (after quality control, z -scoring, and averaging across annotators for a given summary), broken down across the two summarization systems that were used to generate the sample summaries. For CNN/DailyMail, the focus-coverage scores appear to be slightly higher (esp. for BERT), but are better separated over XSUM. We also present the distribution of the inter-sentential coherence scores in Figure 4, and again see that BERT and PG appear to be similar, with PG+Coverage appearing to be slightly better. We return to evaluate these trends more formally in Section 6.3.

²⁷ Emphasizing that there is no task-specific training for any of the BERTSCORE variants; we are simply selecting which layer of which pre-trained model to extract the word representations from.

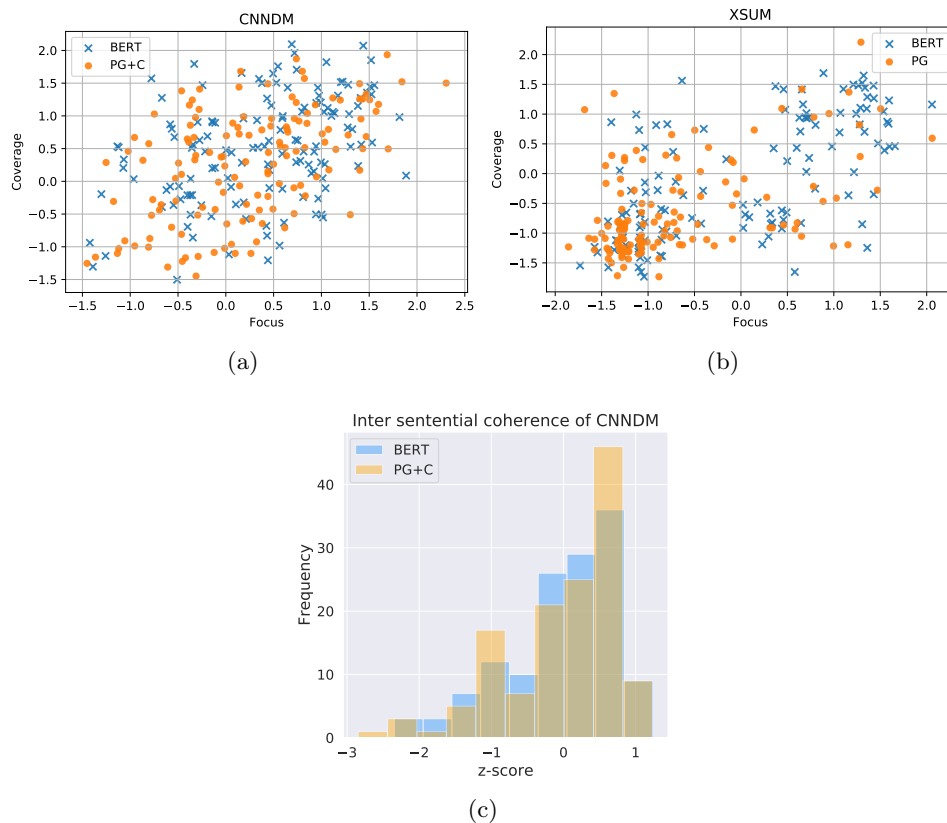


Figure 4: Human annotation result for focus, coverage, and inter-sentential coherence.

6.2.2 METRIC EVALUATION

In Table 8, we present the meta-evaluation results for the primary metrics over focus, coverage, and inter-sentential coherence. We measure Pearson’s r between human judgements and the various automatic metrics discussed in Section 3.

First, we observe that ROUGE, METEOR, and BLEU perform worse than the model-based metrics in all cases.²⁸ For inter-sentential coherence in particular, these baseline metrics perform expectedly badly, around random. Our second observation is that QAGS performs poorly for focus and coverage, compared to traditional metrics.²⁹ The correlation for QAGS is comparable to ROUGE-1 and only slightly better than ROUGE-2. For STS-SCORE, the best segment granularity is sentence, although there is little difference between the three granularities. The results for STS-SCORE are excellent for BERT over XSUM, but appreciably worse for other data-model pairs.

Our optimized version of BERTSCORE performs better than the original due to the task-specific layer selection. Similar to faithfulness (Section 6.1), layer selection is conducted by

28. ROUGE and METEOR scores are calculated based on the original implementations, while BLEU is based on SacreBLEU (Post, 2018).

29. We experimented with different numbers of questions K , ranging from 10 to 50, and also with different methods for pruning ill-formed questions.

Metric	Focus				Coverage				IC	
	C-PG	C-BT	X-PG	X-BT	C-PG	C-BT	X-PG	X-BT	C-PG	C-BT
ROUGE-1	0.607	0.623	0.540	0.562	0.592	0.641	0.480	0.514	0.097	0.138
ROUGE-2	0.595	0.552	0.564	0.454	0.547	0.569	0.463	0.437	-0.004	0.083
ROUGE-LCS	0.604	0.619	0.528	0.552	0.581	0.636	0.482	0.487	0.088	0.114
METEOR	—	—	—	—	0.597	0.660	0.523	0.601	0.061	0.143
BLEU-4	0.511	0.442	0.526	0.304	—	—	—	—	-0.030	0.090
QAGS	0.543	0.611	0.541	0.527	0.570	0.608	0.452	0.513	—	—
STS-SCORE (EDU)	0.525	0.591	0.317	0.527	0.559	0.551	0.461	0.551	—	—
STS-SCORE (sentence)	0.524	0.526	0.444	0.617	0.559	0.572	0.559	0.641	—	—
STS-SCORE (doc)	0.524	0.569	0.444	0.617	0.513	0.508	0.559	0.641	0.124	0.197
BERTSCORE *	0.552	0.519	0.427	0.406	0.549	0.579	0.363	0.359	0.042	0.152
BERTSCORE (Ours)	0.665	0.625	0.577	0.581	0.680	0.695	0.617	0.623	0.055	0.132
Nayeem and Chali (2017)	—	—	—	—	—	—	—	—	-0.275	0.166
NSP	—	—	—	—	—	—	—	—	0.388	0.351

Table 8: Pearson correlation for focus, coverage, and inter-sentential coherence, measured between human judgement and various automatic metrics. (“C-PG” = CNN/DailyMail-PG; “C-BT” = CNN/DailyMail-BERT; “X-PG” = XSUM-PG; “X-BT” = XSUM-BERT; “IC” = Inter-sentential Coherence). All focus metrics are precision based, coverage metrics are recall based, and baselines for IC use F1. ‘*’ uses `roberta-large` (layer 24), while ours use `gpt2-xl` (focus: layer 29, coverage: layer 4, IC: layer 47). Nayeem and Chali (2017) use $\lambda = 0.5$.

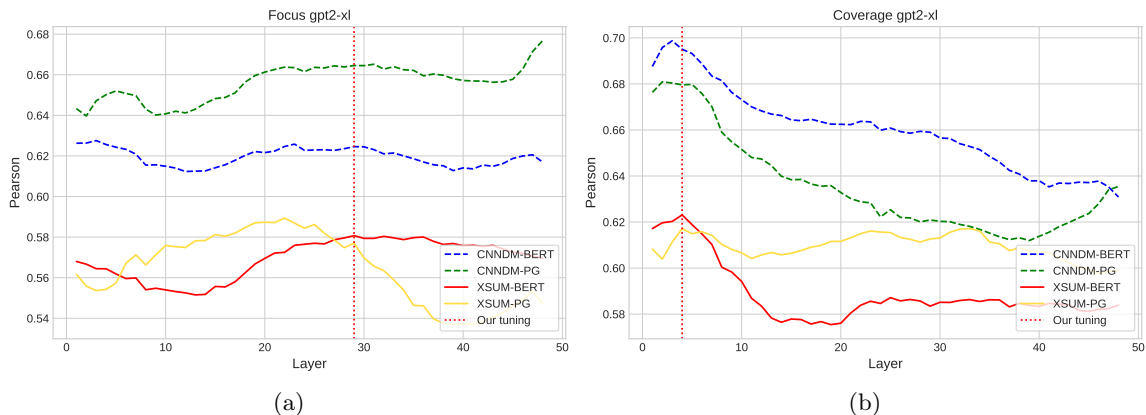


Figure 5: Pearson correlation for each layer of `gpt2-xl` for focus and coverage evaluation. Tuning refers to layer selection (i.e. model parameters are not updated.) Zhang et al. (2020b) does not report WMT tuning for this model.

selecting the average rank of $2 \text{ data} \times 2 \text{ summarization models} = 4$ options, for each 19 pre-trained models. We found that our BERTSCORE (`gpt2-xl`) performs the best for focus

Method	ROUGE			FFCI			
	R-1	R-2	R-L	Fa	Fo	C	IC
LEAD3	40.1	17.3	36.3	91.2	49.2	70.9	65.3
Abstractive							
PG (See et al., 2017)	36.4	15.7	33.4	90.9	52.1	65.6	52.8
PG+C (See et al., 2017)	39.5	17.3	36.4	91.1	52.4	68.6	67.2
rnn+RL+rerank (Chen & Bansal, 2018)	40.9	17.8	38.5	89.6	53.4	70.2	56.4
BOTTOM-UP (Gehrmann et al., 2018)	41.5	18.7	38.6	90.0	55.3	68.5	65.3
BERTSUMEXTABS (Liu & Lapata, 2019)	42.1	19.4	39.1	89.8	51.9	68.7	65.7
BART (Lewis et al., 2020)	44.3	21.1	41.2	89.5	52.6	69.5	69.6
PEGASUS (Zhang et al., 2020a)	44.4	21.5	41.4	89.9	56.0	70.8	69.5
PROPHETNET (Yan et al., 2020)	44.4	21.2	41.5	89.9	55.9	72.0	70.0
Extractive							
BanditSum (Dong et al., 2018)	41.6	18.7	37.9	91.8	51.5	71.6	61.5
PNBERT (Zhong et al., 2019)	42.7	19.5	38.8	91.9	51.9	73.5	66.2
BERTSUMEXT (Liu & Lapata, 2019)	43.3	20.2	39.7	91.8	52.2	73.0	61.8
MATCHSUM (Zhong et al., 2020)	44.4	20.8	40.6	91.9	53.3	72.4	62.5

Table 9: ROUGE and FFCI scores for various summarization models over CNN/DailyMail (“Fa” = faithfulness; “Fo” = focus; “C” = coverage; and “IC” = inter-sentential coherence).

(layer 29) and coverage (layer 4) as depicted in Figure 5. We refer readers who are interested in other pre-trained language model scores to the Appendix .

Finally, we show the effectiveness of NSP prediction for inter-sentential coherence. Computing BERTSCORE against the reference results in low correlation, around random for CNN/DailyMail-PG. We found that a simple NSP score consistently outperforms coherence score (Nayeem & Chali, 2017) at nearly double the correlation score (for CNN/DailyMail-BERT). For CNN/DailyMail-PG, Nayeem and Chali (2017) produces negative correlation $r = -0.28$ which we suspect is due to severe repetition in PG summaries, and the influence of NESim in the coherence score equation. Our proposed NSP score handles this case better and achieves $r = 0.388$ for CNN/DailyMail-PG.

6.3 Summarization Leaderboard Using FFCI

Having motivated the FFCI framework and developed robust metrics for each of the four elements, we next apply them in evaluating a broad range of contemporary methods over CNN/DailyMail and XSUM.

First, we collect summaries from the different abstractive and extractive summarization models either by downloading test outputs provided by the authors or applying checkpoint models from the authors to the test data to generate summaries. In each case, we ensure the test summaries result in similar ROUGE scores to those reported by the authors.

Method	ROUGE			FFCI			
	R-1	R-2	R-L	Fa	Fo	C	IC
LEAD1	16.3	1.6	12.0	90.3	35.3	50.1	—
PG (See et al., 2017)	29.7	9.2	23.2	85.2	45.0	57.1	—
TCONV (Narayan et al., 2018b)	31.9	11.5	25.8	85.2	49.4	57.7	—
BERTSUMEXTABS (Liu & Lapata, 2019)	38.8	16.5	31.3	85.6	53.7	62.3	—
BART (Lewis et al., 2020)	45.1	22.3	37.3	86.6	61.9	69.0	—
PEGASUS (Zhang et al., 2020a)	47.2	24.6	39.3	86.5	64.6	69.5	—

Table 10: ROUGE and FFCI scores for various summarization models over XSUM (“Fa” = faithfulness; “Fo” = focus; “C” = coverage; and “IC” = inter-sentential coherence).

In Tables 9 and 10 we provide FFCI-based results over CNN/DailyMail and XSUM respectively. At a glance, we can see that, despite the lacklustre results for ROUGE in our meta-evaluation, model development based on ROUGE has broadly led to positive progress in summarization, but we get a richer picture of the relative advantages of different methods.

First, the upper bound of faithfulness in CNN/DailyMail is around 91.0, as indicated by LEAD3 and the extractive models. Most abstractive models except PG+C obtain a faithfulness score lower than 91.0, but none lower than 89.0. As such, for CNN/DailyMail, faithfulness appears not to be a differentiating factor, although there has been a slight downward creep with recent abstractive methods. For XSUM, the upper bound for faithfulness is 90.3, for LEAD1. We observe the faithfulness gap between LEAD1 and neural models is bigger than CNN/DailyMail, at around 5–6 points, but there has been a very slight upward trend in faithfulness for abstractive models.

In terms of coverage for abstractive methods, for CNN/DailyMail there has been little improvement in recent years, with all models achieving below the LEAD3 baseline until the recently-proposed PROPHETNET (Yan et al., 2020). Where progress has occurred for abstractive models over CNN/DailyMail is in focus, although results have fluctuated, with BERTSUMEXTABS (Liu & Lapata, 2019) performing notably badly in terms of focus, and to a lesser degree, coverage. This is despite ROUGE suggesting that the model performs better than BOTTOM-UP (Gehrmann et al., 2018), for example. Another example of a substantial change-up in results is BART vs. PEGASUS, where our FFCI framework shows that PEGASUS is substantially better in terms of both focus and coverage, despite the ROUGE scores being almost identical.

In contrast with the abstractive models, the extractive models tend to have higher coverage and lower focus. While MATCHSUM (Zhong et al., 2020) is state of the art in terms of ROUGE, based on our evaluation, coverage is actually markedly lower than competitor methods but focus is high.

For XSUM, we can observe large improvements in focus in particular, and relatively smaller but still clear improvements in coverage. Faithfulness, on the other hand, has improved only slightly, and there is clear room for improvement. Once again, our framework

Model	IC Score		
	Automatic	Human (<i>z</i> -score)	
		Sum	Avg
LEAD3	65.3	0.45	-0.0030
PG+C	67.2	-3.40	-0.0026
BERTSUMEXTABS	65.7	-4.50	-0.0330
PROPHETNET	70.0	7.50	0.0600

Table 11: Inter-sentential coherence (IC) based on automatic and manual scores.

clearly shows that PEGASUS achieves better focus than BART, but that they are closer in terms of coverage.

Lastly, we look at inter-sentential coherence for CNN/DailyMail.³⁰ Overall, three abstractive models achieve IC score close to LEAD3: PG+C, BOTTOM-UP, and BERTSUMEXTABS. PG and rnn+RL+rerank (Chen & Bansal, 2018) result in very poor inter-sentential coherence, which we suspect is due to severe repetition at the decoding stage. BART, PEGASUS and PROPHETNET achieve higher scores than LEAD3, with a gap of around 4 points. As expected, the extractive methods tend to result in poorer inter-sentential coherence than LEAD3, with the exception of PNBERT.

6.4 Analysis of Inter-Sentential Coherence

One surprising observation from Table 9 is that the IC score of LEAD3 is actually lower than that of PROPHETNET, despite LEAD3 consisting of the first three sentences of the source document (which we would assume have high IC). Additionally, the raw correlation numbers for the IC experiments from Table 8 are low, suggesting that the scores from the IC metrics are prone to noise and potentially unreliable. Given this, we performed additional manual annotations for IC over four summarization methods to better understand the results: LEAD3, PG+C, BERTSUMEXTABS, and PROPHETNET. We use the same 135 CNN/DailyMail samples as in Section 5.1, and ask three workers to annotate IC using the same procedure as Section 5.2 (135 documents \times 4 models \times 3 annotators, resulting in 1,620 annotations). In this annotation, we achieve a quality score of 96.7 and average working time of 31.7 minutes.

Firstly, in Table 11 we observe that the manual annotations for PG+C and BERTSUMEXTABS are consistent with Figure 4 from Section 6.2, in that PG+C tends to have slightly higher inter-sentential coherence. More importantly, Table 11 confirms that PROPHETNET does indeed outperform all other models including LEAD3. The reason that LEAD3 has lower IC is that the resultant summaries often contain subtle disfluencies, as shown in Table 12. The first LEAD3 example contains a “teaser” first sentence (which is disconnected from the next two sentences), while the second example has metadata in the first sentence. Compared with this, the PROPHETNET examples are more fluent in terms of information structure.

30. Recalling that XSUM summaries are single sentence, and thus inter-sentential coherence is irrelevant.

Lead3	ProphetNet
<ul style="list-style-type: none"> • This is the breathtaking moment a diver came face to face with a 'fish tornado' off the mexican coast. • Tori hester, 25, from san diego, california, was diving in cabo pulmo when the huge school of trevally fish began circling above her. • Husband jeff, 26, was on hand to capture the incredible moment using his underwater camera. 	<ul style="list-style-type: none"> • Tori hester, 25, from san diego, california, was diving in cabo pulmo. • Huge school of trevally fish began circling above her. • Husband jeff, a marine scientist, was on hand to capture the moment.
<ul style="list-style-type: none"> • London (CNN). • Congolese immigrant tarsis mboma thale has a small business selling t-shirts in johannesburg, south africa. • Thale's job normally requires him to walk the streets of the city he has called home for the past few years. 	<ul style="list-style-type: none"> • Wave of anti-immigrant violence has swept south africa in recent days, leaving several dead. • Some blame alleged inflammatory comments about foreign nationals from the zulu king. • Others say a labor dispute between locals and foreigners back in march turned nasty.

Table 12: Example LEAD3 and PROPHETNET summaries for CNN/DailyMail.

Although the analysis in Table 11 is highly promising in terms of the veracity of the automatic metric, this is the least-developed of the four FFCI metrics with plenty of room for further improvement (owing to its low absolute correlation values (0.35–0.38), as seen in Table 8).

7. Conclusion

We introduce the FFCI evaluation framework for summarization evaluation, based on the four elements of: faithfulness, focus, coverage, and inter-sentential coherence. We have shown that BERTSCORE (roberta-base) is the most robust metric for evaluating faithfulness, BERTSCORE (gpt2-x1) for focus and coverage, and NSP-score for inter-sentential coherence.

Our general finding is that ROUGE has lead to positive progress in modern summarization systems but lacks fine-grained interpretability. FFCI shows that since the LSTM-based seq2seq, modern abstractive summarization systems over CNN/DailyMail have largely improved on focus, with coverage not being much better than LEAD3 until recent systems (e.g. PROPHETNET). Our FFCI framework found three competitive state-of-the-art systems: BART, PEGASUS, and PROPHETNET, with PEGASUS and PROPHETNET having generally higher focus and coverage respectively.

Lastly, although FFCI was designed based on our survey of evaluation approaches in previous work (Table 1), we believe there are some additional aspects that should be addressed in future work such as redundancy and relevance. Our work and the aforementioned survey

are mostly based on news datasets such as CNN/DM and XSUM with relatively short summaries. We believe the redundancy in particular becomes very important as summarisation research shifts focus to longer documents and summaries.

Acknowledgments

In this research, the first author is supported by the Australia Awards Scholarship (AAS), funded by the Department of Foreign Affairs and Trade (DFAT), Australia. This research was undertaken using the LIEF HPC-GPGPU Facility hosted at The University of Melbourne. This facility was established with the assistance of LIEF Grant LE170100200.

Appendix A. Recommended Layers for Faithfulness, Focus, and Coverage

Model	Fa	Fo	C
bert-base-uncased	6	1	2
bert-large-uncased	11	9	9
roberta-base	10	9	2
roberta-large	13	13	3
roberta-large-mnli	14	15	3
xlnet-base-cased	6	4	2
xlnet-large-cased	7	7	5
gpt2	1	3	3
gpt2-medium	8	5	1
gpt2-large	2	21	3
gpt2-xl	2	29	4
t5-small	2	3	2
t5-base	3	4	4
t5-large	10	13	10
bart-base	1	3	1
bart-large	2	5	2
pegasus-xsum	8	11	6
pegasus-cnn_dailymail	12	11	5
pegasus-large	3	4	4

Table 13: Recommended layers for faithfulness (Fa), focus (Fo), and coverage (Co).

As discussed by Zhang et al. (2020b) and Reimers and Gurevych (2019b), layer selection in BERT model is important. BERTSCORE was designed to maximize the Pearson correlation between $\mathcal{F}_{\text{BERT}}$ and WMT16, which is potentially less than optimal for evaluating focus and coverage in summarization.

In Table 13, we present 19 models and their recommended layer number for evaluating faithfulness, focus, and coverage. Specifically, we use FA_{MODEL} , $\mathcal{P}_{\text{MODEL}}$ and $\mathcal{R}_{\text{MODEL}}$ to calculate the Pearson correlation, respectively. To pick the best layer, we simply average the rank of each data-model based on the outputs of a given layer. We observe that almost all selected layers are different to BERTSCORE (see Figures 6–12 for examples). The optimal layer for focus tends to be one of the last layers, while earlier layers tend to work better for coverage. We also present the layer selection plots for the non-BERT models (Figures 13–24)

Appendix B. Pre-Trained Language Model Scores of Faithfulness

Model	Faithfulness			
	PG	TRANS2S	TCONV	BERT
<i>Results based on our layer selection</i>				
bert-base-uncased	0.424	0.394	0.460	0.463
bert-large-uncased	0.420	0.406	0.436	0.473
roberta-base	0.459	0.450	0.519	0.475
roberta-large	0.411	0.425	0.474	0.489
roberta-large-mnli	0.437	0.415	0.489	0.477
xlnet-base-cased	0.355	0.347	0.372	0.373
xlnet-large-cased	0.369	0.378	0.393	0.386
gpt2	0.299	0.341	0.331	0.367
gpt2-medium	0.329	0.412	0.357	0.396
gpt2-large	0.360	0.399	0.386	0.439
gpt2-xl	0.357	0.392	0.381	0.431
t5-small	0.326	0.307	0.330	0.328
t5-base	0.334	0.332	0.346	0.353
t5-large	0.346	0.344	0.355	0.354
bart-base	0.370	0.383	0.381	0.421
bart-large	0.375	0.412	0.405	0.452
pegasus-xsum	0.406	0.410	0.437	0.417
pegasus-cnn.dailymail	0.401	0.406	0.432	0.413
pegasus-large	0.392	0.417	0.387	0.463
<i>Results based on recommended layers by Zhang et al. (2020b)</i>				
bert-base-uncased	0.386	0.377	0.435	0.440
bert-large-uncased	0.251	0.335	0.380	0.378
roberta-base	0.459	0.450	0.519	0.475
roberta-large	0.150	0.168	0.162	0.230
roberta-large-mnli	0.370	0.340	0.440	0.423
xlnet-base-cased	0.281	0.303	0.355	0.328
xlnet-large-cased	0.369	0.378	0.393	0.386

Table 14: Pearson correlation of all experimental results on pre-trained language model scores for faithfulness (XSUM data). We highlight models with the highest average across data-model pairs. **roberta-base** of Zhang et al. (2020b) and ours use the same layer-10. Please note that Zhang et al. (2020b) final recommendation is to use **roberta-large** (layer-24).

Appendix C. Pre-Trained Language Model Scores of Focus and Coverage

Model	Focus				Coverage			
	C-PG	C-BT	X-PG	X-BT	C-PG	C-BT	X-PG	X-BT
<i>Results based on layer selection</i>								
bert-base-uncased	0.623	0.647	0.513	0.529	0.636	0.680	0.587	0.578
bert-large-uncased	0.627	0.641	0.547	0.563	0.648	0.689	0.608	0.609
roberta-base	0.621	0.636	0.531	0.550	0.698	0.707	0.553	0.596
roberta-large	0.634	0.643	0.583	0.552	0.674	0.712	0.588	0.603
roberta-large-mnli	0.647	0.658	0.573	0.557	0.677	0.706	0.591	0.610
xlnet-base-cased	0.640	0.603	0.508	0.531	0.636	0.639	0.566	0.533
xlnet-large-cased	0.636	0.612	0.522	0.552	0.638	0.651	0.585	0.554
gpt2	0.621	0.620	0.493	0.528	0.648	0.678	0.534	0.562
gpt2-medium	0.636	0.616	0.579	0.544	0.667	0.687	0.552	0.603
gpt2-large	0.668	0.629	0.577	0.571	0.676	0.689	0.614	0.612
gpt2-xl	0.665	0.625	0.577	0.581	0.680	0.695	0.617	0.623
t5-small	0.632	0.603	0.529	0.582	0.641	0.636	0.580	0.591
t5-base	0.632	0.608	0.548	0.586	0.641	0.653	0.580	0.596
t5-large	0.643	0.615	0.544	0.591	0.646	0.641	0.600	0.604
bart-base	0.664	0.629	0.541	0.546	0.674	0.688	0.577	0.578
bart-large	0.655	0.627	0.550	0.561	0.676	0.690	0.599	0.609
pegasus-xsum	0.630	0.643	0.556	0.563	0.664	0.676	0.587	0.603
pegasus-cnn_dailymail	0.652	0.640	0.534	0.571	0.668	0.700	0.580	0.604
pegasus-large	0.625	0.626	0.561	0.563	0.669	0.712	0.597	0.624
<i>Results based on recommended layers by Zhang et al. (2020b)</i>								
bert-base-uncased	0.613	0.624	0.478	0.532	0.583	0.665	0.601	0.553
bert-large-uncased	0.602	0.630	0.502	0.524	0.610	0.648	0.619	0.524
roberta-base	0.616	0.638	0.516	0.556	0.677	0.667	0.544	0.559
roberta-large	0.619	0.641	0.584	0.548	0.694	0.691	0.562	0.577
roberta-large-mnli	0.608	0.653	0.529	0.537	0.628	0.668	0.562	0.554
xlnet-base-cased	0.644	0.608	0.461	0.507	0.603	0.628	0.529	0.489
xlnet-large-cased	0.636	0.612	0.522	0.552	0.633	0.639	0.584	0.558

Table 15: Pearson correlation of pre-trained language model scores for focus and coverage (“C-PG” = CNNDM-PG; “C-BT” = CNNDM-BERT; “X-PG” = XSUM-PG; “X-BT” = XSUM-BERT). We highlight models with the highest average across data-model pairs.

Appendix D. Full Results Over Inter-Sentential Coherence

Model	PG	BERT
NSP-Score (mean)		
bert-base-uncased	0.388 ± 0.069	0.351 ± 0.051
roberta-base	0.339 ± 0.037	0.230 ± 0.061
albert-base-v2	0.331 ± 0.049	0.200 ± 0.045
xlnet-base-cased	0.365 ± 0.051	0.235 ± 0.070
electra-base-discriminator	0.389 ± 0.053	0.305 ± 0.038
gpt2	0.313 ± 0.008	0.114 ± 0.024
bart-base	0.357 ± 0.069	0.256 ± 0.068
NSP-Score (max)		
bert-base-uncased	0.269 ± 0.085	0.342 ± 0.045
roberta-base	0.339 ± 0.037	0.230 ± 0.061
albert-base-v2	0.336 ± 0.052	0.243 ± 0.054
xlnet-base-cased	0.247 ± 0.059	0.219 ± 0.069
electra-base-discriminator	0.338 ± 0.032	0.311 ± 0.050
gpt2	0.212 ± 0.013	0.064 ± 0.046
bart-base	0.292 ± 0.055	0.267 ± 0.097
NSP-Score (min)		
bert-base-uncased	0.375 ± 0.049	0.245 ± 0.052
roberta-base	0.261 ± 0.053	0.148 ± 0.075
albert-base-v2	0.293 ± 0.064	0.151 ± 0.032
xlnet-base-cased	0.349 ± 0.046	0.136 ± 0.064
electra-base-discriminator	0.306 ± 0.060	0.227 ± 0.039
gpt2	0.356 ± 0.018	0.142 ± 0.021
bart-base	0.317 ± 0.067	0.171 ± 0.037
Nayeem and Chali (2017)		
$\lambda = 0$	0.046	0.131
$\lambda = 0.3$	-0.193	0.160
$\lambda = 0.5$	-0.275	0.166
$\lambda = 0.7$	-0.312	0.156
$\lambda = 1.0$	-0.334	0.128

Table 16: Pearson correlation of all experimental results on inter-sentential coherence. NSP-Score is computed 5 times over data variant-5 (see Section 5.3).

Appendix E. Pre-Trained Language Model Scores in Different Layers Over Faithfulness, Focus, and Coverage

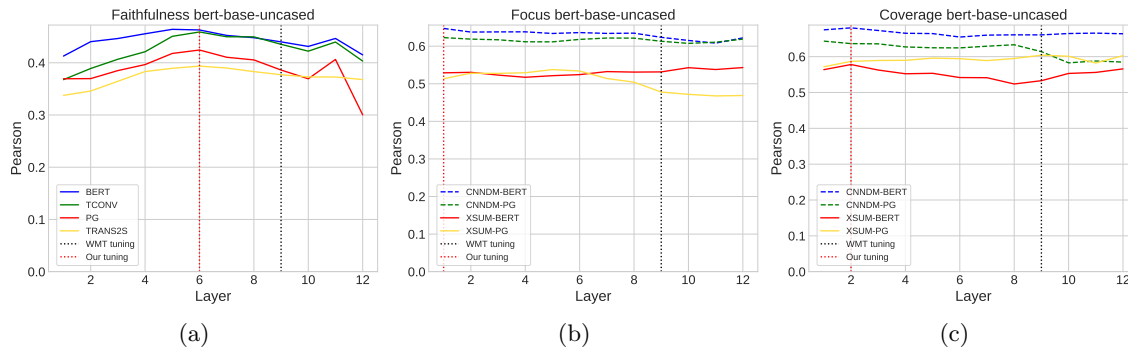


Figure 6: bert-base-uncased

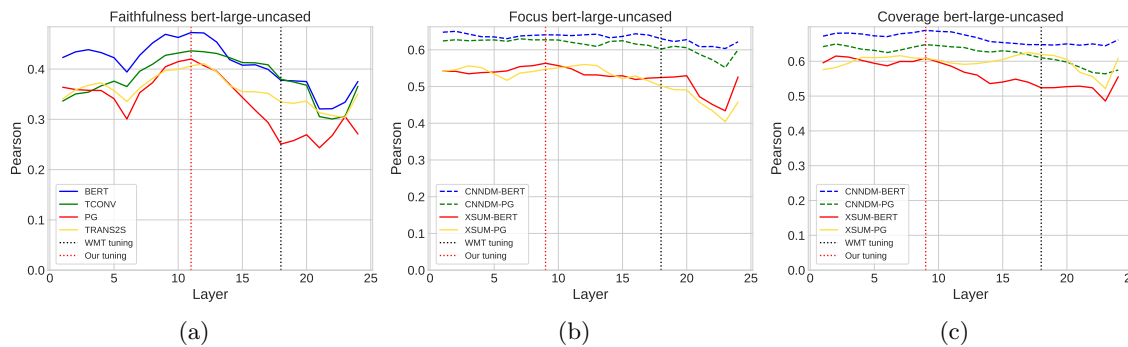


Figure 7: bert-large-uncased

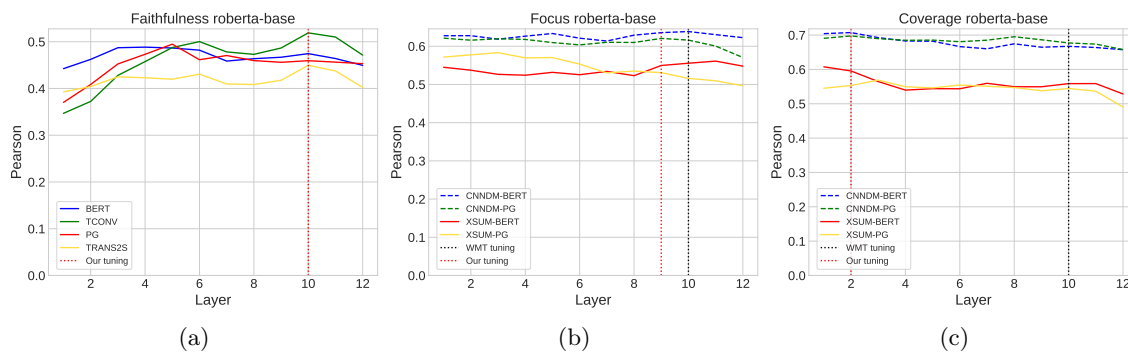


Figure 8: roberta-base

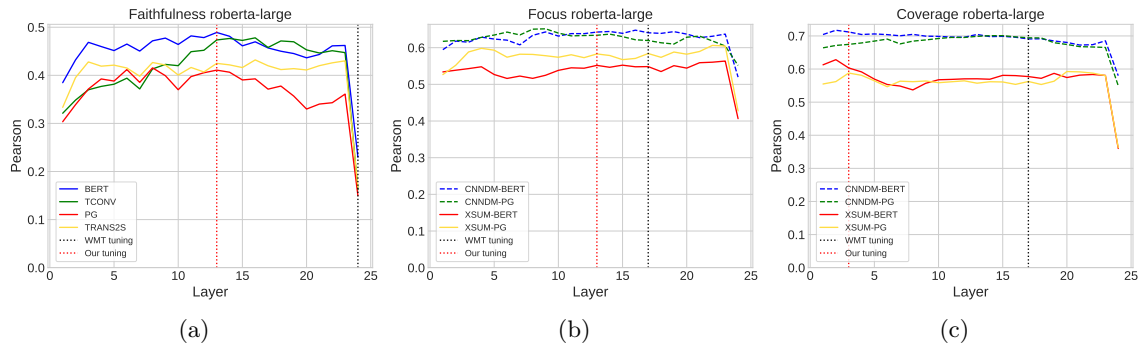


Figure 9: `roberta-large`. Please note that final recommendation of Zhang et al. (2020b) is to use layer-24, however, in their supplementary material (Appendix B), the best layer is 17.

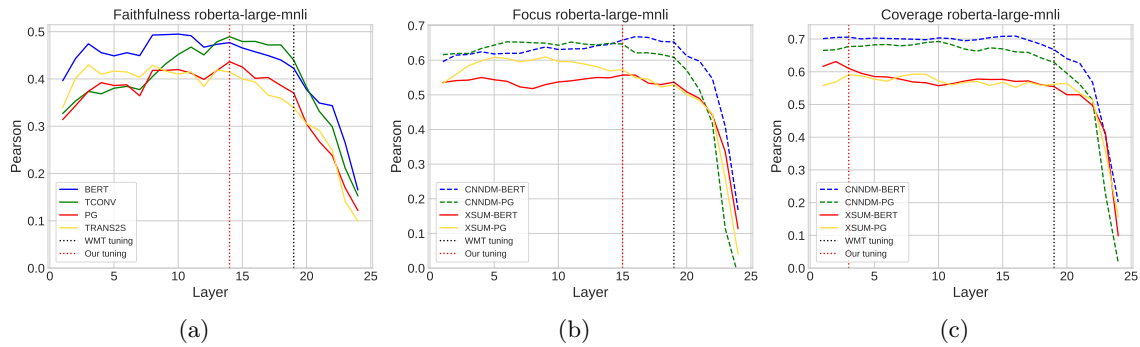


Figure 10: `roberta-large-mnli`.

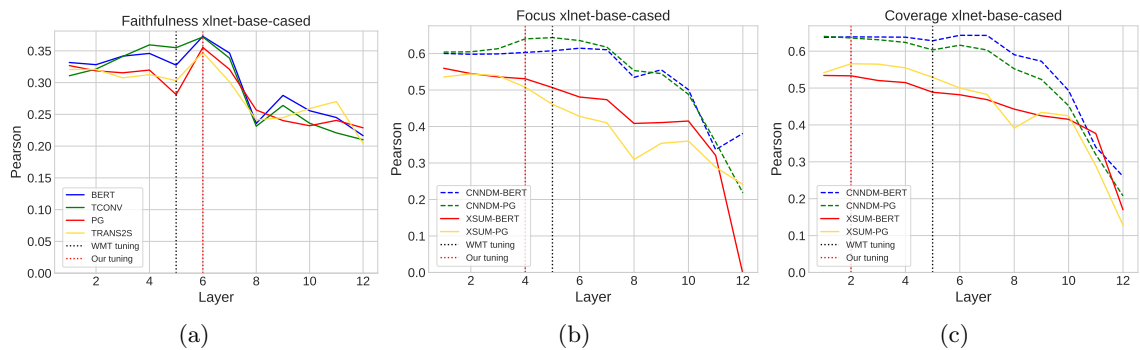


Figure 11: `xlnet-base-cased`.

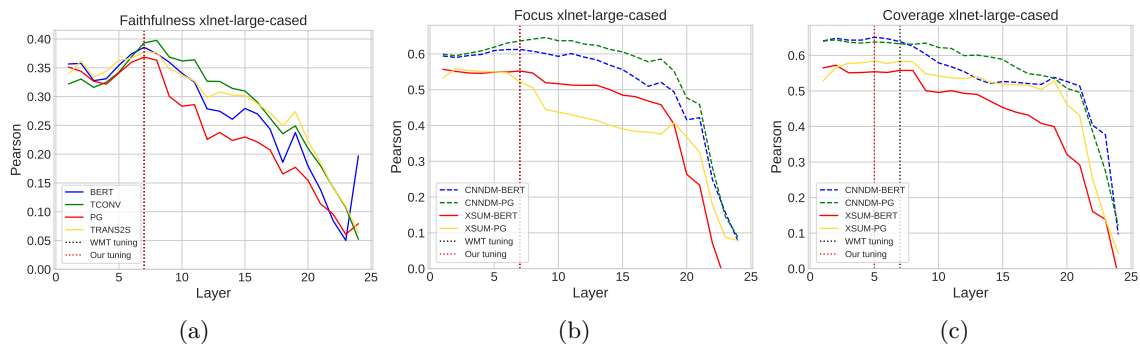


Figure 12: xlnet-large-cased.

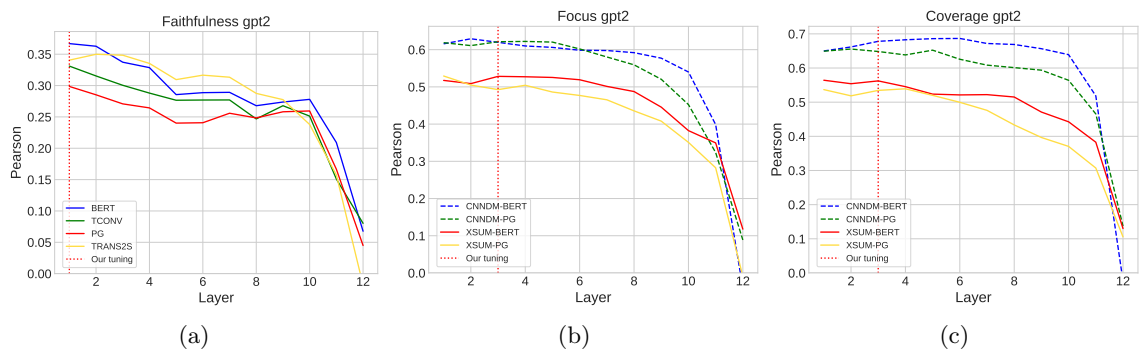


Figure 13: gpt2.

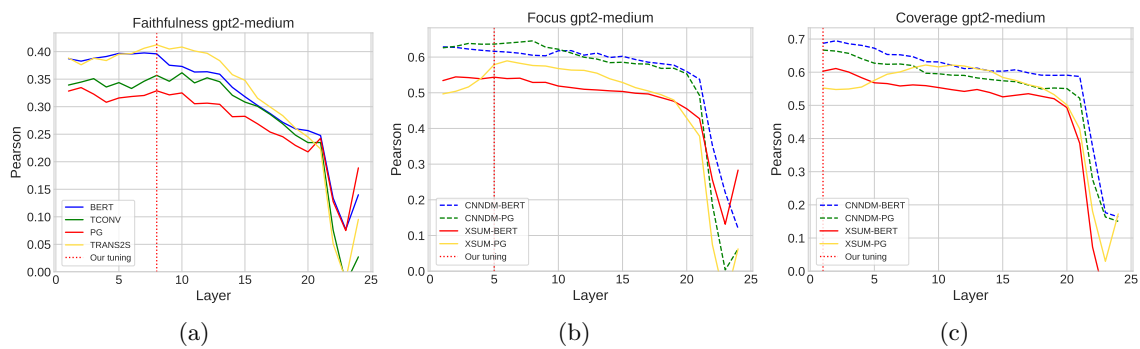


Figure 14: gpt2-medium.

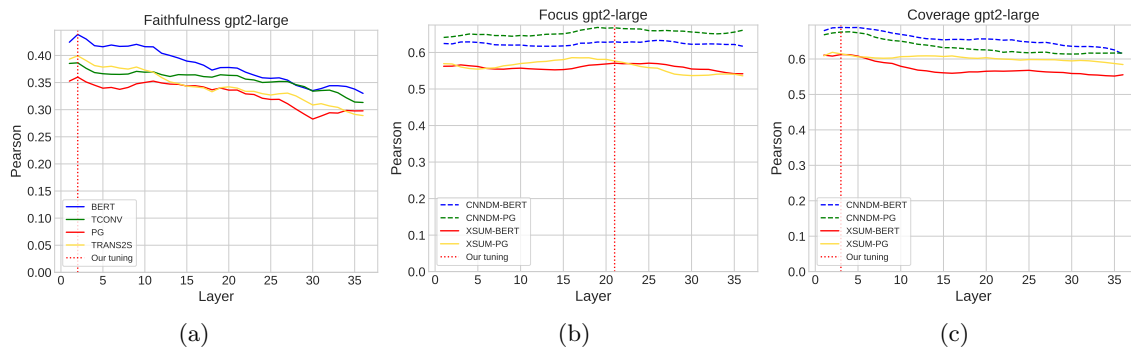


Figure 15: gpt2-large.

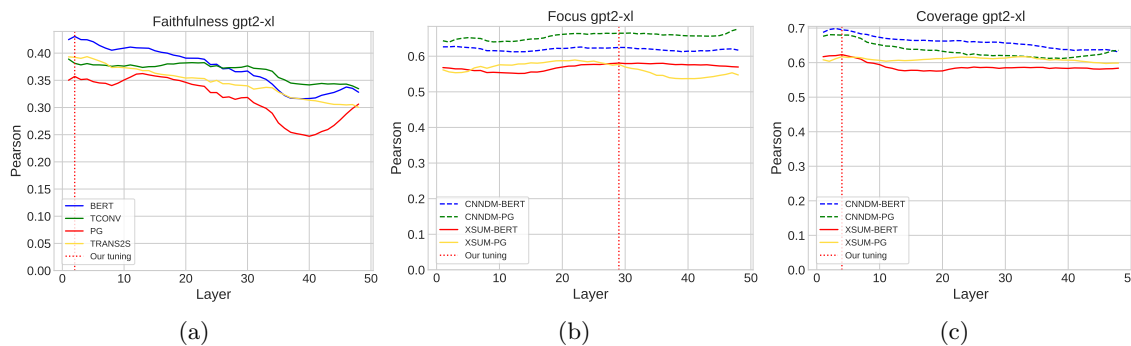


Figure 16: gpt2-xl.

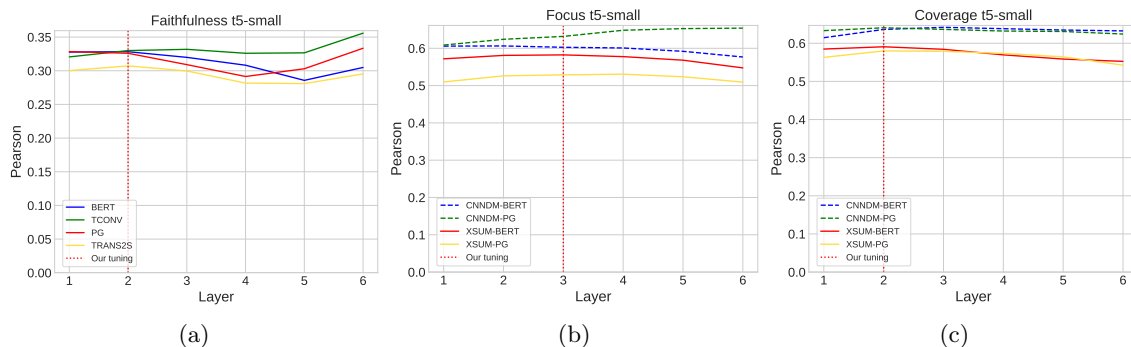


Figure 17: t5-small.

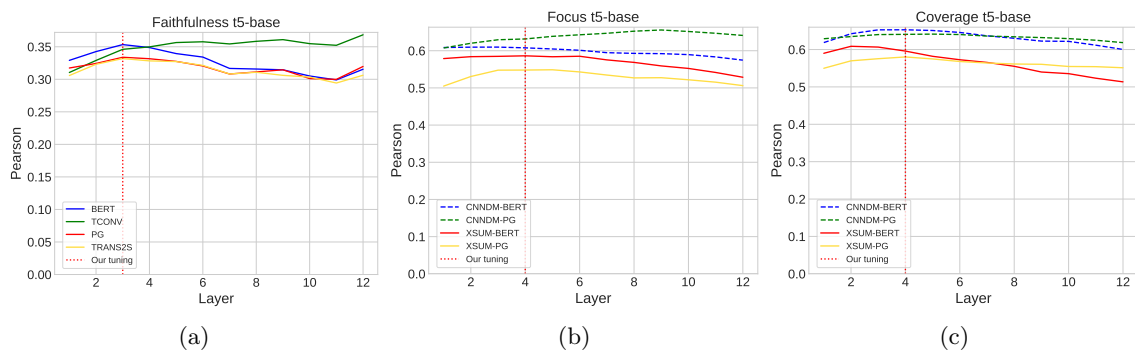


Figure 18: t5-base.

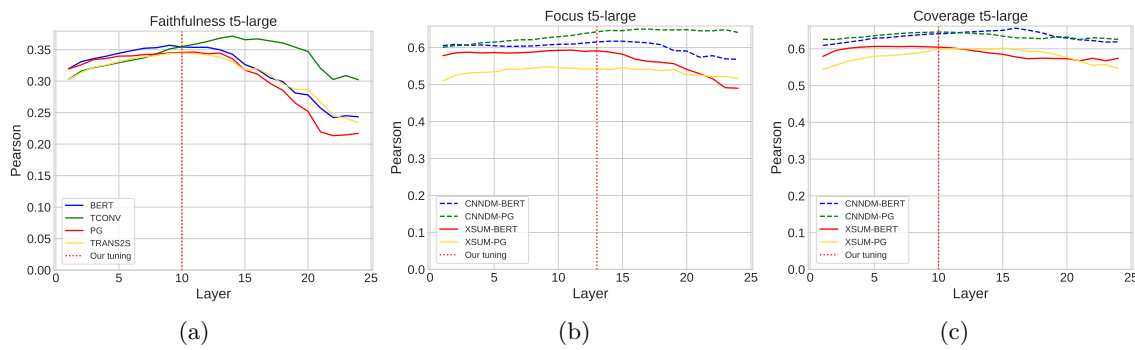


Figure 19: t5-large.

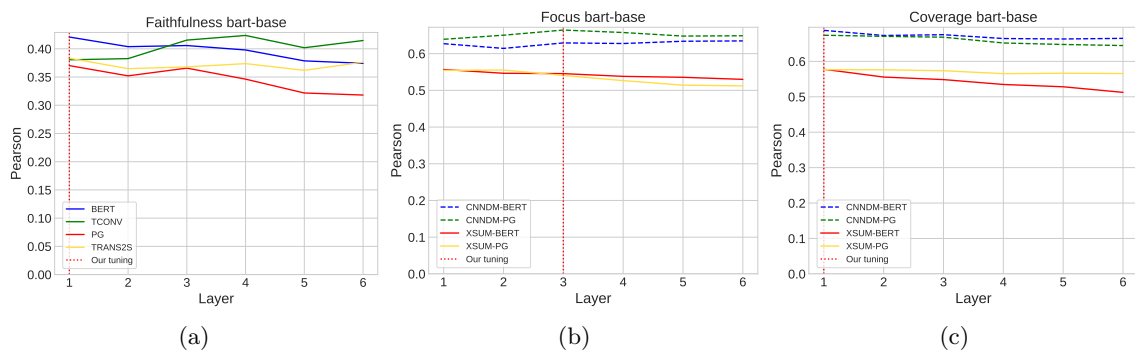


Figure 20: bart-base.

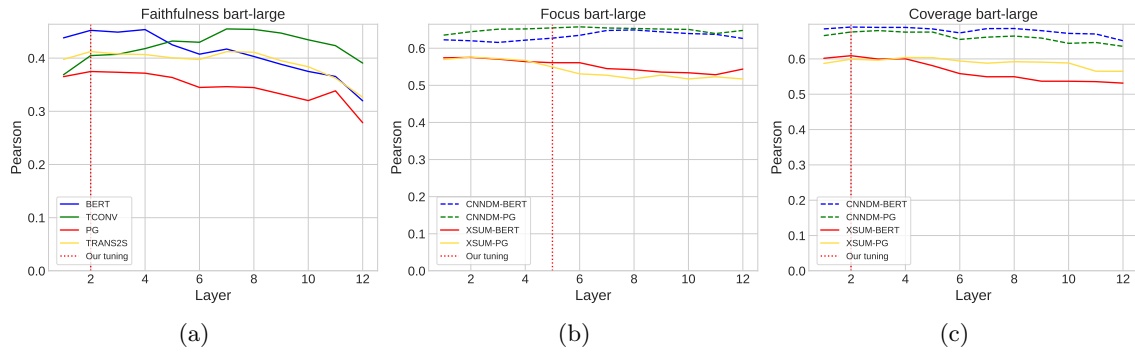


Figure 21: bart-large.

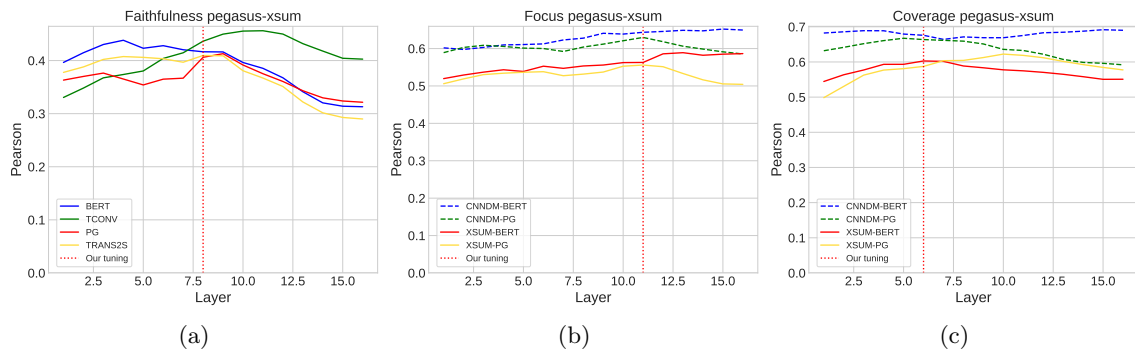


Figure 22: pegasus-xsum.

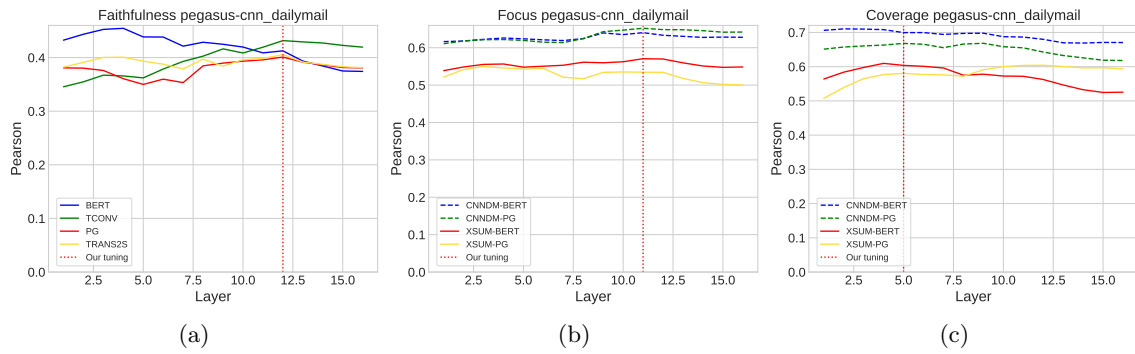


Figure 23: pegasus-cnn.dailymail.

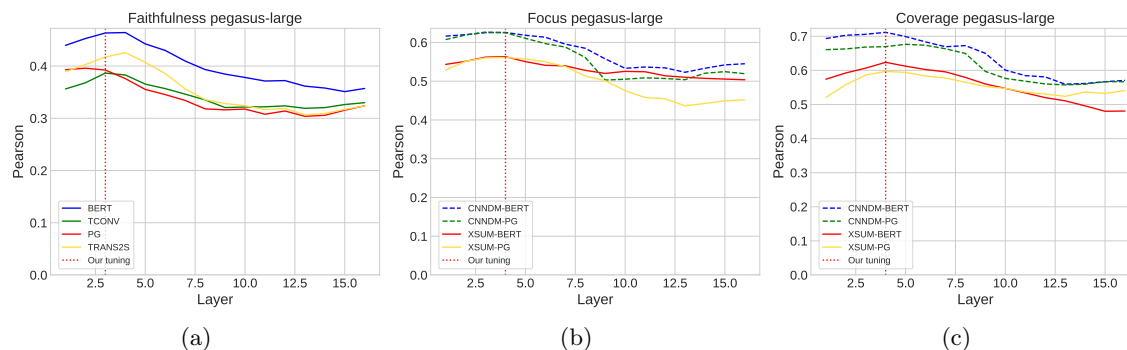


Figure 24: pegasus-large.

References

- Abacha, A. B., & Demner-Fushman, D. (2019). On the summarization of consumer health questions. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2228–2234.
- Agirre, E., Cer, D., Diab, M., & Gonzalez-Agirre, A. (2012). SemEval-2012 task 6: A pilot on semantic textual similarity. In **SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pp. 385–393, Montréal, Canada. Association for Computational Linguistics.
- Ahmad, W., Chakraborty, S., Ray, B., & Chang, K.-W. (2020). A transformer-based approach for source code summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4998–5007, Online. Association for Computational Linguistics.
- Amplayo, R. K., & Lapata, M. (2020). Unsupervised opinion summarization with noising and denoising. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 1934–1945, Online. Association for Computational Linguistics.
- Amplayo, R. K., Lim, S., & Hwang, S.-w. (2018). Entity commonsense representation for neural abstractive summarization. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 697–707, New Orleans, Louisiana. Association for Computational Linguistics.
- Anderson, P., Fernando, B., Johnson, M., & Gould, S. (2016). SPICE: Semantic propositional image caption evaluation. In *European Conference on Computer Vision (14th : 2016)*, pp. 382–398.
- Bhandari, M., Gour, P. N., Ashfaq, A., Liu, P., & Neubig, G. (2020). Re-evaluating evaluation in text summarization. In *Proceedings of the 2020 Conference on Empirical*

- Methods in Natural Language Processing (EMNLP)*, pp. 9347–9359, Online. Association for Computational Linguistics.
- Bommasani, R., & Cardie, C. (2020). Intrinsic evaluation of summarization datasets. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 8075–8096, Online. Association for Computational Linguistics.
- Bražinskas, A., Lapata, M., & Titov, I. (2020). Unsupervised opinion summarization as copycat-review generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 5151–5169, Online. Association for Computational Linguistics.
- Cao, M., Dong, Y., Wu, J., & Cheung, J. C. K. (2020). Factual error correction for abstractive summarization models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6251–6258, Online. Association for Computational Linguistics.
- Cao, Z., Li, W., Li, S., & Wei, F. (2018). Retrieve, rerank and rewrite: Soft template based neural summarization. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 152–161, Melbourne, Australia. Association for Computational Linguistics.
- Celikyilmaz, A., Bosselut, A., He, X., & Choi, Y. (2018). Deep communicating agents for abstractive summarization. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 1662–1675, New Orleans, Louisiana. Association for Computational Linguistics.
- Chen, Y.-C., & Bansal, M. (2018). Fast abstractive summarization with reinforce-selected sentence rewriting. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 675–686, Melbourne, Australia. Association for Computational Linguistics.
- Cohan, A., Deroncourt, F., Kim, D. S., Bui, T., Kim, S., Chang, W., & Goharian, N. (2018). A discourse-aware attention model for abstractive summarization of long documents. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pp. 615–621, New Orleans, Louisiana. Association for Computational Linguistics.
- Cui, P., Hu, L., & Liu, Y. (2020). Enhancing extractive text summarization with topic-aware graph neural networks. In *Proceedings of the 28th International Conference on Computational Linguistics*, pp. 5360–5371, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Deng, Y., Zhang, W., & Lam, W. (2020). Multi-hop inference for question-driven summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6734–6744, Online. Association for Computational Linguistics.
- Desai, S., Xu, J., & Durrett, G. (2020). Compressive summarization with plausibility and salience modeling. In *Proceedings of the 2020 Conference on Empirical Methods in*

- Natural Language Processing (EMNLP)*, pp. 6259–6274, Online. Association for Computational Linguistics.
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Dong, Y., Shen, Y., Crawford, E., van Hoof, H., & Cheung, J. C. K. (2018). Banditsum: Extractive summarization as a contextual bandit. In *EMNLP 2018: 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 3739–3748.
- Du, X., Shao, J., & Cardie, C. (2017). Learning to ask: Neural question generation for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1342–1352, Vancouver, Canada. Association for Computational Linguistics.
- Duan, X., Yin, M., Zhang, M., Chen, B., & Luo, W. (2019a). Zero-shot cross-lingual abstractive sentence summarization through teaching generation and attention. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 3162–3172, Florence, Italy. Association for Computational Linguistics.
- Duan, X., Yu, H., Yin, M., Zhang, M., Luo, W., & Zhang, Y. (2019b). Contrastive attention mechanism for abstractive sentence summarization. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3044–3053, Hong Kong, China. Association for Computational Linguistics.
- Durmus, E., He, H., & Diab, M. (2020). FEQA: A question answering evaluation framework for faithfulness assessment in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 5055–5070, Online. Association for Computational Linguistics.
- Fabbri, A., Li, I., She, T., Li, S., & Radev, D. (2019). Multi-news: A large-scale multi-document summarization dataset and abstractive hierarchical model. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 1074–1084, Florence, Italy. Association for Computational Linguistics.
- Fabbri, A. R., Kryściński, W., McCann, B., Xiong, C., Socher, R., & Radev, D. (2020). Summeval: Re-evaluating summarization evaluation. *arXiv preprint arXiv:2007.12626*.
- Falke, T., & Gurevych, I. (2019). Fast concept mention grouping for concept map-based multi-document summarization. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 695–700, Minneapolis, Minnesota. Association for Computational Linguistics.
- Frermann, L., & Klementiev, A. (2019). Inducing document structure for aspect-based summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 6263–6273, Florence, Italy. Association for Computational Linguistics.

- Gao, S., Chen, X., Li, P., Chan, Z., Zhao, D., & Yan, R. (2019). How to write summaries with patterns? learning towards abstractive summarization through prototype editing. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3741–3751, Hong Kong, China. Association for Computational Linguistics.
- Gao, Y., Zhao, W., & Eger, S. (2020). SUPERT: Towards new frontiers in unsupervised evaluation metrics for multi-document summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 1347–1354, Online. Association for Computational Linguistics.
- Gehrmann, S., Deng, Y., & Rush, A. (2018). Bottom-up abstractive summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 4098–4109, Brussels, Belgium. Association for Computational Linguistics.
- Gholipour Ghalandari, D., Hokamp, C., Pham, N. T., Glover, J., & Ifrim, G. (2020). A large-scale multi-document summarization dataset from the Wikipedia current events portal. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 1302–1308, Online. Association for Computational Linguistics.
- Gholipour Ghalandari, D., & Ifrim, G. (2020). Examining the state-of-the-art in news timeline summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 1322–1334, Online. Association for Computational Linguistics.
- Graham, Y., Baldwin, T., Dowling, M., Eskevich, M., Lynn, T., & Tounsi, L. (2016). Is all that glitters in machine translation quality estimation really gold?. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pp. 3124–3134, Osaka, Japan. The COLING 2016 Organizing Committee.
- Graham, Y., Baldwin, T., & Mathur, N. (2015). Accurate evaluation of segment-level machine translation metrics. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1183–1191, Denver, Colorado. Association for Computational Linguistics.
- Graham, Y., Baldwin, T., Moffat, A., & Zobel, J. (2017). Can machine translation systems be evaluated by the crowd alone. *Natural Language Engineering*, 23(1), 3–30.
- Grenander, M., Dong, Y., Cheung, J. C. K., & Louis, A. (2019). Countering the effects of lead bias in news summarization via multi-stage training and auxiliary losses. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 6019–6024, Hong Kong, China. Association for Computational Linguistics.
- Gui, M., Tian, J., Wang, R., & Yang, Z. (2019). Attention optimization for abstractive document summarization. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on*

- Natural Language Processing (EMNLP-IJCNLP)*, pp. 1222–1228, Hong Kong, China. Association for Computational Linguistics.
- Guo, H., Pasunuru, R., & Bansal, M. (2018). Soft layer-specific multi-task summarization with entailment and question generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 687–697, Melbourne, Australia. Association for Computational Linguistics.
- Hardy, H., Narayan, S., & Vlachos, A. (2019). HighRES: Highlight-based reference-less evaluation of summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 3381–3392, Florence, Italy. Association for Computational Linguistics.
- Hardy, H., & Vlachos, A. (2018). Guided neural language generation for abstractive summarization using Abstract Meaning Representation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 768–773, Brussels, Belgium. Association for Computational Linguistics.
- He, R., Zhao, L., & Liu, H. (2020). TWEETSUM: Event oriented social summarization dataset. In *Proceedings of the 28th International Conference on Computational Linguistics*, pp. 5731–5736, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Hermann, K. M., Kocisky, T., Grefenstette, E., Espeholt, L., Kay, W., Suleyman, M., & Blunsom, P. (2015). Teaching machines to read and comprehend. In *Neural Information Processing Systems*, pp. 1693–1701.
- Hsu, W.-T., Lin, C.-K., Lee, M.-Y., Min, K., Tang, J., & Sun, M. (2018). A unified model for extractive and abstractive summarization using inconsistency loss. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 132–141, Melbourne, Australia. Association for Computational Linguistics.
- Huang, D., Cui, L., Yang, S., Bao, G., Wang, K., Xie, J., & Zhang, Y. (2020a). What have we achieved on text summarization?. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 446–469, Online. Association for Computational Linguistics.
- Huang, K.-H., Li, C., & Chang, K.-W. (2020b). Generating sports news from live commentary: A Chinese dataset for sports game summarization. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pp. 609–615, Suzhou, China. Association for Computational Linguistics.
- Huang, L., Wu, L., & Wang, L. (2020c). Knowledge graph-augmented abstractive summarization with semantic-driven cloze reward. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 5094–5107, Online. Association for Computational Linguistics.
- Isonuma, M., Mori, J., & Sakata, I. (2019). Unsupervised neural single-document summarization of reviews via learning latent discourse structure and its ranking. In *Proceed-*

- ings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2142–2152, Florence, Italy. Association for Computational Linguistics.
- Jadhav, A., & Rajan, V. (2018). Extractive summarization with SWAP-NET: Sentences and words from alternating pointer networks. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 142–151, Melbourne, Australia. Association for Computational Linguistics.
- Ji, Y., & Eisenstein, J. (2014). Representation learning for text-level discourse parsing. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 13–24, Baltimore, Maryland. Association for Computational Linguistics.
- Jia, R., Rajpurkar, P., & Liang, P. (2018). Know what you don't know: Unanswerable questions for squad. In *ACL 2018: 56th Annual Meeting of the Association for Computational Linguistics*, Vol. 2, pp. 784–789.
- Jia, R., Cao, Y., Tang, H., Fang, F., Cao, C., & Wang, S. (2020). Neural extractive summarization with hierarchical attentive heterogeneous graph network. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 3622–3631, Online. Association for Computational Linguistics.
- Kano, R., Miura, Y., Taniguchi, T., & Ohkuma, T. (2020). Identifying implicit quotes for unsupervised extractive summarization of conversations. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pp. 291–302, Suzhou, China. Association for Computational Linguistics.
- Kedzie, C., McKeown, K., & Daumé III, H. (2018). Content selection in deep learning models of summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 1818–1828, Brussels, Belgium. Association for Computational Linguistics.
- Kim, B., Kim, H., & Kim, G. (2019). Abstractive summarization of Reddit posts with multi-level memory networks. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 2519–2531, Minneapolis, Minnesota. Association for Computational Linguistics.
- Koto, F., Lau, J. H., & Baldwin, T. (2020). Liputan6: A large-scale Indonesian dataset for text summarization. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pp. 598–608, Suzhou, China. Association for Computational Linguistics.
- Kouris, P., Alexandridis, G., & Stafylopatis, A. (2019). Abstractive text summarization based on deep learning and semantic content generalization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 5082–5092, Florence, Italy. Association for Computational Linguistics.
- Krishna, K., & Srinivasan, B. V. (2018). Generating topic-oriented summaries using neural attention. In *Proceedings of the 2018 Conference of the North American Chapter of the*

Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pp. 1697–1705, New Orleans, Louisiana. Association for Computational Linguistics.

- Kryscinski, W., McCann, B., Xiong, C., & Socher, R. (2020). Evaluating the factual consistency of abstractive text summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 9332–9346, Online. Association for Computational Linguistics.
- Kryściński, W., Paulus, R., Xiong, C., & Socher, R. (2018). Improving abstraction in text summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 1808–1817, Brussels, Belgium. Association for Computational Linguistics.
- Ladhak, F., Li, B., Al-Onaizan, Y., & McKeown, K. (2020). Exploring content selection in summarization of novel chapters. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 5043–5054, Online. Association for Computational Linguistics.
- Lavie, A., & Agarwal, A. (2007). METEOR: An automatic metric for MT evaluation with high levels of correlation with human judgments. In *Proceedings of the Second Workshop on Statistical Machine Translation*, pp. 228–231, Prague, Czech Republic. Association for Computational Linguistics.
- Lebanoff, L., Deroncourt, F., Kim, D. S., Chang, W., & Liu, F. (2020). A cascade approach to neural abstractive summarization with content selection and fusion. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pp. 529–535, Suzhou, China. Association for Computational Linguistics.
- Lebanoff, L., Song, K., Deroncourt, F., Kim, D. S., Kim, S., Chang, W., & Liu, F. (2019). Scoring sentence singletons and pairs for abstractive summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2175–2189, Florence, Italy. Association for Computational Linguistics.
- Lee, D., Shin, M. C., Whang, T., Cho, S., Ko, B., Lee, D., Kim, E., & Jo, J. (2020). Reference and document aware semantic evaluation methods for Korean language summarization. In *Proceedings of the 28th International Conference on Computational Linguistics*, pp. 5604–5616, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Lev, G., Shmueli-Scheuer, M., Herzig, J., Jerbi, A., & Konopnicki, D. (2019). TalkSumm: A dataset and scalable annotation method for scientific paper summarization based on conference talks. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2125–2131, Florence, Italy. Association for Computational Linguistics.
- Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Stoyanov, V., & Zettlemoyer, L. (2020). BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 7871–7880.

- Li, C., Xu, W., Li, S., & Gao, S. (2018a). Guiding generation for abstractive text summarization based on key information guide network. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pp. 55–60, New Orleans, Louisiana. Association for Computational Linguistics.
- Li, H., Zhu, J., Zhang, J., & Zong, C. (2018b). Ensure the correctness of the summary: Incorporate entailment knowledge into abstractive sentence summarization. In *Proceedings of the 27th International Conference on Computational Linguistics*, pp. 1430–1441, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Li, M., Zhang, L., Ji, H., & Radke, R. J. (2019a). Keep meeting summaries on topic: Abstractive multi-modal meeting summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2190–2196, Florence, Italy. Association for Computational Linguistics.
- Li, S., Lei, D., Qin, P., & Wang, W. Y. (2019b). Deep reinforcement learning with distributional semantic rewards for abstractive summarization. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 6038–6044, Hong Kong, China. Association for Computational Linguistics.
- Li, Z., Wu, W., & Li, S. (2020). Composing elementary discourse units in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 6191–6196, Online. Association for Computational Linguistics.
- Liao, K., Lebanoff, L., & Liu, F. (2018). Abstract Meaning Representation for multi-document summarization. In *Proceedings of the 27th International Conference on Computational Linguistics*, pp. 1178–1190, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Lin, C.-Y. (2004). ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pp. 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Lin, J., Sun, X., Ma, S., & Su, Q. (2018). Global encoding for abstractive summarization. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 163–169, Melbourne, Australia. Association for Computational Linguistics.
- Liu, Y., & Lapata, M. (2019). Text summarization with pretrained encoders. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3730–3740, Hong Kong, China. Association for Computational Linguistics.
- Liu, Y., Titov, I., & Lapata, M. (2019a). Single document summarization as tree induction. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 1745–1755, Minneapolis, Minnesota. Association for Computational Linguistics.

- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., & Stoyanov, V. (2019b). RoBERTa: A robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Lu, Y., Dong, Y., & Charlin, L. (2020). Multi-XScience: A large-scale dataset for extreme multi-document summarization of scientific articles. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 8068–8074, Online. Association for Computational Linguistics.
- Luo, L., Ao, X., Song, Y., Pan, F., Yang, M., & He, Q. (2019). Reading like HER: Human reading inspired extractive summarization. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3033–3043, Hong Kong, China. Association for Computational Linguistics.
- Makino, T., Iwakura, T., Takamura, H., & Okumura, M. (2019). Global optimization under length constraint for neural text summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 1039–1048, Florence, Italy. Association for Computational Linguistics.
- Mann, W. C. (1984). Discourse structures for text generation. In *10th International Conference on Computational Linguistics and 22nd Annual Meeting of the Association for Computational Linguistics*, pp. 367–375, Stanford, California, USA. Association for Computational Linguistics.
- Mao, Y., Liu, L., Zhu, Q., Ren, X., & Han, J. (2020a). Facet-aware evaluation for extractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4941–4957, Online. Association for Computational Linguistics.
- Mao, Y., Qu, Y., Xie, Y., Ren, X., & Han, J. (2020b). Multi-document summarization with maximal marginal relevance-guided reinforcement learning. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1737–1751, Online. Association for Computational Linguistics.
- Mathur, N., Baldwin, T., & Cohn, T. (2019). Putting evaluation in context: Contextual embeddings improve machine translation evaluation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2799–2808, Florence, Italy. Association for Computational Linguistics.
- Maynez, J., Narayan, S., Bohnet, B., & McDonald, R. (2020). On faithfulness and factuality in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 1906–1919, Online. Association for Computational Linguistics.
- Mendes, A., Narayan, S., Miranda, S., Marinho, Z., Martins, A. F. T., & Cohen, S. B. (2019). Jointly extracting and compressing documents with summary state representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 3955–3966, Minneapolis, Minnesota. Association for Computational Linguistics.

- Mikolov, T., Yih, W.-t., & Zweig, G. (2013). Linguistic regularities in continuous space word representations. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 746–751, Atlanta, Georgia. Association for Computational Linguistics.
- Nallapati, R., Zhou, B., dos Santos, C., Gulçehre, Ç., & Xiang, B. (2016). Abstractive text summarization using sequence-to-sequence RNNs and beyond. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, pp. 280–290, Berlin, Germany. Association for Computational Linguistics.
- Narayan, S., Cardenas, R., Papasarantopoulos, N., Cohen, S. B., Lapata, M., Yu, J., & Chang, Y. (2018a). Document modeling with external attention for sentence extraction. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2020–2030, Melbourne, Australia. Association for Computational Linguistics.
- Narayan, S., Cohen, S. B., & Lapata, M. (2018b). Don’t give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 1797–1807, Brussels, Belgium. Association for Computational Linguistics.
- Narayan, S., Cohen, S. B., & Lapata, M. (2018c). Ranking sentences for extractive summarization with reinforcement learning. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 1747–1759, New Orleans, Louisiana. Association for Computational Linguistics.
- Nayeem, M. T., & Chali, Y. (2017). Extract with order for coherent multi-document summarization. In *Proceedings of TextGraphs-11: the Workshop on Graph-based Methods for Natural Language Processing*, pp. 51–56, Vancouver, Canada. Association for Computational Linguistics.
- Nenkova, A., & Passonneau, R. (2004). Evaluating content selection in summarization: The pyramid method. In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pp. 145–152, Boston, Massachusetts, USA. Association for Computational Linguistics.
- Ng, J.-P., & Abrecht, V. (2015). Better summarization evaluation with word embeddings for ROUGE. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pp. 1925–1930, Lisbon, Portugal. Association for Computational Linguistics.
- Ouyang, J., Song, B., & McKeown, K. (2019). A robust abstractive system for cross-lingual summarization. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 2025–2031, Minneapolis, Minnesota. Association for Computational Linguistics.
- Pagnoni, A., Balachandran, V., & Tsvetkov, Y. (2021). Understanding factuality in abstractive summarization with FRANK: A benchmark for factuality metrics. In *Proceedings of the 2021 Conference of the North American Chapter of the Association*

- for *Computational Linguistics: Human Language Technologies*, pp. 4812–4829, Online. Association for Computational Linguistics.
- Palaskar, S., Libovický, J., Gella, S., & Metze, F. (2019). Multimodal abstractive summarization for how2 videos. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 6587–6596, Florence, Italy. Association for Computational Linguistics.
- Papalampidi, P., Keller, F., Frermann, L., & Lapata, M. (2020). Screenplay summarization using latent narrative structure. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 1920–1933, Online. Association for Computational Linguistics.
- Papalampidi, P., Keller, F., & Lapata, M. (2019). Movie plot analysis via turning point identification. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 1707–1717, Hong Kong, China. Association for Computational Linguistics.
- Papineni, K., Roukos, S., Ward, T., & Zhu, W.-J. (2002). Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pp. 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Parida, S., & Motlicek, P. (2019). Abstract text summarization: A low resource challenge. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 5994–5998, Hong Kong, China. Association for Computational Linguistics.
- Pasunuru, R., & Bansal, M. (2018). Multi-reward reinforced summarization with saliency and entailment. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pp. 646–653, New Orleans, Louisiana. Association for Computational Linguistics.
- Pennington, J., Socher, R., & Manning, C. (2014). GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Peyrard, M. (2019a). A simple theoretical model of importance for summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 1059–1073, Florence, Italy. Association for Computational Linguistics.
- Peyrard, M. (2019b). Studying summarization evaluation metrics in the appropriate scoring range. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 5093–5100, Florence, Italy. Association for Computational Linguistics.
- Peyrard, M., & Gurevych, I. (2018). Objective function learning to match human judgements for optimization-based summarization. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics:*

- Human Language Technologies, Volume 2 (Short Papers)*, pp. 654–660, New Orleans, Louisiana. Association for Computational Linguistics.
- Pilault, J., Li, R., Subramanian, S., & Pal, C. (2020). On extractive and abstractive neural document summarization with transformer language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 9308–9319, Online. Association for Computational Linguistics.
- Post, M. (2018). A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pp. 186–191, Brussels, Belgium. Association for Computational Linguistics.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8), 9.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., & Liu, P. J. (2019). Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140), 1–67.
- Reimers, N., & Gurevych, I. (2019a). Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Reimers, N., & Gurevych, I. (2019b). Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.
- Sakaue, S., Hirao, T., Nishino, M., & Nagata, M. (2018). Provable fast greedy compressive summarization with any monotone submodular function. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 1737–1746, New Orleans, Louisiana. Association for Computational Linguistics.
- Schumann, R., Mou, L., Lu, Y., Vechtomova, O., & Markert, K. (2020). Discrete optimization for unsupervised sentence summarization with word-level extraction. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 5032–5042, Online. Association for Computational Linguistics.
- Scialom, T., Dray, P.-A., Lamprier, S., Piwowarski, B., & Staiano, J. (2020). MLSUM: The multilingual summarization corpus. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 8051–8067, Online. Association for Computational Linguistics.
- See, A., Liu, P. J., & Manning, C. D. (2017). Get to the point: Summarization with pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1073–1083, Vancouver, Canada. Association for Computational Linguistics.
- ShafieiBavani, E., Ebrahimi, M., Wong, R., & Chen, F. (2018). Summarization evaluation in the absence of human model summaries using the compositionality of word embeddings. In *Proceedings of the 27th International Conference on Computational*

- Linguistics*, pp. 905–914, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Shapira, O., Gabay, D., Gao, Y., Ronen, H., Pasunuru, R., Bansal, M., Amsterdamer, Y., & Dagan, I. (2019). Crowdsourcing lightweight pyramids for manual summary evaluation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 682–687, Minneapolis, Minnesota. Association for Computational Linguistics.
- Sharma, E., Li, C., & Wang, L. (2019). BIGPATENT: A large-scale dataset for abstractive and coherent summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2204–2213, Florence, Italy. Association for Computational Linguistics.
- Shen, A., & Baldwin, T. (2021). A simple yet effective method for sentence ordering. In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pp. 154–160.
- Shen, X., Zhao, Y., Su, H., & Klakow, D. (2019). Improving latent alignment in text summarization by generalizing the pointer generator. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3762–3773, Hong Kong, China. Association for Computational Linguistics.
- Song, K., Zhao, L., & Liu, F. (2018). Structure-infused copy mechanisms for abstractive summarization. In *Proceedings of the 27th International Conference on Computational Linguistics*, pp. 1717–1729, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Sotudeh Gharebagh, S., Goharian, N., & Filice, R. (2020). Attend to medical ontologies: Content selection for clinical abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 1899–1905, Online. Association for Computational Linguistics.
- Suhara, Y., Wang, X., Angelidis, S., & Tan, W.-C. (2020). OpinionDigest: A simple framework for opinion summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 5789–5798, Online. Association for Computational Linguistics.
- Sun, S., & Nenkova, A. (2019). The feasibility of embedding based automatic evaluation for single document summarization. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 1216–1221, Hong Kong, China. Association for Computational Linguistics.
- Tan, B., Qin, L., Xing, E., & Hu, Z. (2020). Summarizing text on any aspects: A knowledge-informed weakly-supervised approach. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6301–6309, Online. Association for Computational Linguistics.

- Trischler, A., Wang, T., Yuan, X., Harris, J., Sordoni, A., Bachman, P., & Suleman, K. (2017). Newsqa: A machine comprehension dataset. In *Proceedings of the 2nd Workshop on Representation Learning for NLP*, pp. 191–200.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2017). Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pp. 5998–6008.
- Vedantam, R., Zitnick, C. L., & Parikh, D. (2015). CIDEr: Consensus-based image description evaluation. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4566–4575.
- Wang, A., Cho, K., & Lewis, M. (2020). Asking and answering questions to evaluate the factual consistency of summaries. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 5008–5020. Association for Computational Linguistics.
- Wang, H., Wang, X., Xiong, W., Yu, M., Guo, X., Chang, S., & Wang, W. Y. (2019a). Self-supervised learning for contextualized extractive summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2221–2227, Florence, Italy. Association for Computational Linguistics.
- Wang, K., Quan, X., & Wang, R. (2019b). BiSET: Bi-directional selective encoding with template for abstractive summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2153–2162, Florence, Italy. Association for Computational Linguistics.
- Wang, K., Chang, B., & Sui, Z. (2020a). A spectral method for unsupervised multi-document summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 435–445, Online. Association for Computational Linguistics.
- Wang, Z., Duan, Z., Zhang, H., Wang, C., Tian, L., Chen, B., & Zhou, M. (2020b). Friendly topic assistant for transformer based abstractive summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 485–497, Online. Association for Computational Linguistics.
- West, P., Holtzman, A., Buys, J., & Choi, Y. (2019). BottleSum: Unsupervised and self-supervised sentence summarization using the information bottleneck principle. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3752–3761, Hong Kong, China. Association for Computational Linguistics.
- Wu, H., Ma, T., Wu, L., Manyumwa, T., & Ji, S. (2020). Unsupervised reference-free summary quality evaluation via contrastive learning. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 3612–3621, Online. Association for Computational Linguistics.
- Xiao, L., Wang, L., He, H., & Jin, Y. (2020). Modeling content importance for summarization with pre-trained language models. In *Proceedings of the 2020 Conference on*

- Empirical Methods in Natural Language Processing (EMNLP)*, pp. 3606–3611, Online. Association for Computational Linguistics.
- Xiao, W., & Carenini, G. (2019). Extractive summarization of long documents by combining global and local context. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3011–3021, Hong Kong, China. Association for Computational Linguistics.
- Xiao, W., & Carenini, G. (2020). Systematically exploring redundancy reduction in summarizing long documents. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pp. 516–528, Suzhou, China. Association for Computational Linguistics.
- Xu, J., Gan, Z., Cheng, Y., & Liu, J. (2020). Discourse-aware neural extractive text summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 5021–5031, Online. Association for Computational Linguistics.
- Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., Zemel, R., & Bengio, Y. (2015). Show, attend and tell: Neural image caption generation with visual attention. In *Proceedings of The 32nd International Conference on Machine Learning*, Vol. 3, pp. 2048–2057.
- Xu, R., Zhu, C., Shi, Y., Zeng, M., & Huang, X. (2020a). Mixed-lingual pre-training for cross-lingual summarization. In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pp. 536–541, Suzhou, China. Association for Computational Linguistics.
- Xu, S., Li, H., Yuan, P., Wu, Y., He, X., & Zhou, B. (2020b). Self-attention guided copy mechanism for abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 1355–1362, Online. Association for Computational Linguistics.
- Xu, Y., & Lapata, M. (2020). Coarse-to-fine query focused multi-document summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 3632–3645, Online. Association for Computational Linguistics.
- Yan, Y., Qi, W., Gong, Y., Liu, D., Duan, N., Chen, J., Zhang, R., & Zhou, M. (2020). Prophetnet: Predicting future n-gram for sequence-to-sequence pre-training. *arXiv preprint arXiv:2001.04063*.
- Yang, Y., Bao, F., & Nenkova, A. (2017). Detecting (un)important content for single-document news summarization. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pp. 707–712, Valencia, Spain. Association for Computational Linguistics.
- Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R., & Le, Q. V. (2019). Xlnet: Generalized autoregressive pretraining for language understanding. In *NeurIPS 2019: Thirty-third Conference on Neural Information Processing Systems*, pp. 5753–5763.

- Yin, Y., Meng, F., Su, J., Ge, Y., Song, L., Zhou, J., & Luo, J. (2020). Enhancing pointer network for sentence ordering with pairwise ordering predictions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 9482–9489.
- You, Y., Jia, W., Liu, T., & Yang, W. (2019). Improving abstractive document summarization with salient information modeling. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2132–2141, Florence, Italy. Association for Computational Linguistics.
- Zhang, J., Zhao, Y., Saleh, M., & Liu, P. (2020a). Pegasus: Pre-training with extracted gap-sentences for abstractive summarization. In *ICML 2020: 37th International Conference on Machine Learning*.
- Zhang, T., Kishore, V., Wu, F., Weinberger, K. Q., & Artzi, Y. (2020b). Bertscore: Evaluating text generation with bert. In *ICLR 2020 : Eighth International Conference on Learning Representations*.
- Zhang, X., Wei, F., & Zhou, M. (2019). HIBERT: Document level pre-training of hierarchical bidirectional transformers for document summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 5059–5069, Florence, Italy. Association for Computational Linguistics.
- Zhao, L., Xu, W., & Guo, J. (2020). Improving abstractive dialogue summarization with graph structures and topic words. In *Proceedings of the 28th International Conference on Computational Linguistics*, pp. 437–449, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Zhao, W., Peyrard, M., Liu, F., Gao, Y., Meyer, C. M., & Eger, S. (2019). MoverScore: Text generation evaluating with contextualized embeddings and earth mover distance. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 563–578, Hong Kong, China. Association for Computational Linguistics.
- Zheng, H., & Lapata, M. (2019). Sentence centrality revisited for unsupervised summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 6236–6247, Florence, Italy. Association for Computational Linguistics.
- Zhong, M., Liu, P., Chen, Y., Wang, D., Qiu, X., & Huang, X. (2020). Extractive summarization as text matching. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 6197–6208, Online. Association for Computational Linguistics.
- Zhong, M., Liu, P., Wang, D., Qiu, X., & Huang, X. (2019). Searching for effective neural extractive summarization: What works and what’s next. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 1049–1058, Florence, Italy. Association for Computational Linguistics.
- Zhou, J., & Rush, A. (2019). Simple unsupervised summarization by contextual matching. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 5101–5106, Florence, Italy. Association for Computational Linguistics.

- Zhu, J., Wang, Q., Wang, Y., Zhou, Y., Zhang, J., Wang, S., & Zong, C. (2019). NCLS: Neural cross-lingual summarization. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3054–3064, Hong Kong, China. Association for Computational Linguistics.
- Zhu, J., Zhou, Y., Zhang, J., & Zong, C. (2020). Attend, translate and summarize: An efficient method for neural cross-lingual summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 1309–1321, Online. Association for Computational Linguistics.
- Zou, Y., Zhang, X., Lu, W., Wei, F., & Zhou, M. (2020). Pre-training for abstractive document summarization by reinstating source text. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 3646–3660, Online. Association for Computational Linguistics.